

Hanoi University of Science and Technology  
School of Information and Communication Technology

# **Task Scheduling In Real-time System**

Quang Khanh  
khanh.tq170083@gmail.com

2nd June 2021



System overview

Problem Statement

Solution

Experiment result

# Table of Contents



System overview

Problem Statement

Solution

Experiment result

# Google's Borg System



**Google's Borg system is a cluster manager** that runs millions of jobs, from many thousands of different applications, across a number of clusters each with up to ten thousands of machines

**Borg cell** is a set of machines grouped together

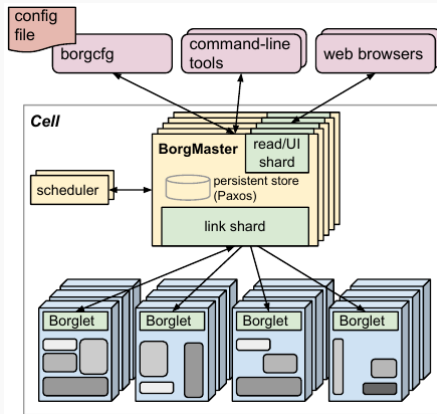


Figure: Borg architecture

# Google's Borg System



**BorgMaster** is a logically centralized controller over a borg cell

**Borglet** is an agent process that runs on each machine in a cell

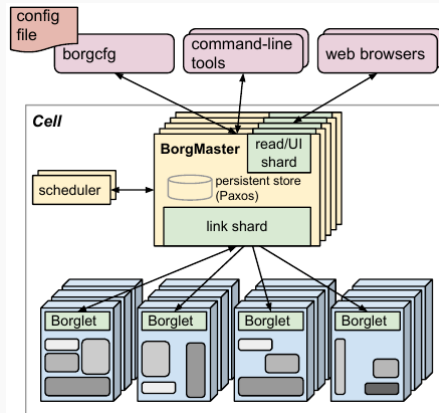


Figure: Borg architecture

# Task scheduling in Borg system



## When a job is submitted

It is added to pending queue

## Scheduler

Scheduler scans the pending queue and assigns tasks to machines by scheduling algorithm. The algorithm has two parts:

- ▶ **feasible checking:** to find machines on which the task could run
- ▶ **scoring:** to pick one machine to assign task to.

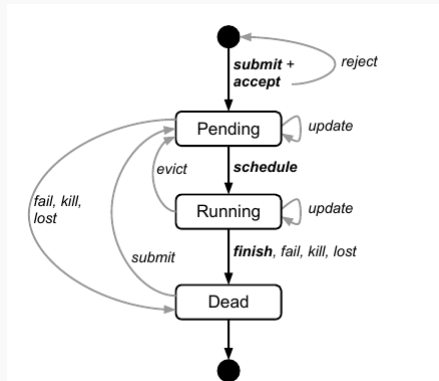


Figure: Borg architecture



## Feasible checking

The scheduler find a set of machines that:

- ▶ meet the task's constraints (dependencies, packages, ...)
- ▶ have enough available resources

## Scoring

The score is based on some criteria such as minimizing the number and priority of preempted tasks, and resources-related criteria are:

- ▶ **worst fit:** speading workload across all the machines uniformly
- ▶ **best fit:** tries to fill machines as tightly as possible

# Table of Contents



System overview

**Problem Statement**

Solution

Experiment result





## Execution time

Due to lack of certainty about how long the tasks' instructions are, the scheduler cannot estimate execution time of the tasks

## Service quantity

The solution cannot optimize execution time, leading to longer waiting-time to user



## Eviction

When resources of machines are not available, some running tasks would be evicted to release resources for higher-priority tasks

## Service quality

The user will have to wait longer to finish their job

## System effectiveness

The resources spent for the evicted tasks would be wasted. The google trace data published in 2011 shows that up to 60% of CPU resource is wasted by tasks that do not complete successfully.

# Table of Contents



System overview

Problem Statement

**Solution**

Experiment result



## Divide tasks to two types:

- ▶ Long-running: the task as an always running service, unspecified its instructions
- ▶ Batch-job (short-running): the task which has specific instructions depending on only its input, which takes from a few seconds to a few days to complete.

## Long-running tasks

Find feasible machines and pick one of them for the task

## Short-running tasks

Predict the intructions and find a machine in which the execution time is minimized

# Proposed task scheduling pipeline

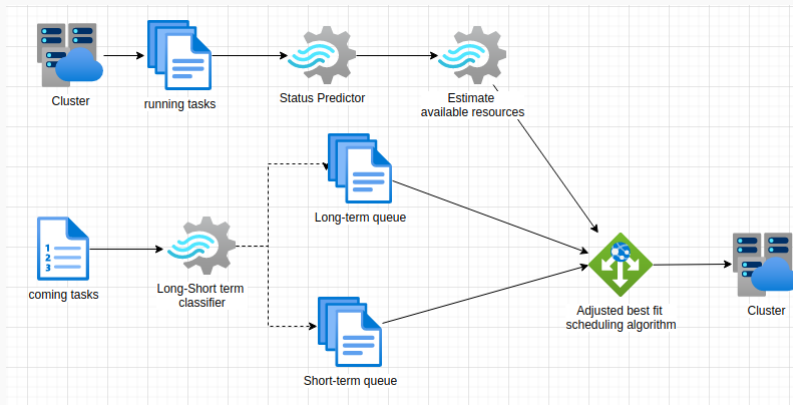


Figure: Task scheduling pipeline

## Prior-knowledge classifier

Basing on what kind of applications submitting tasks

## Logistic Classifier

A logistic regression model to predict label of tasks

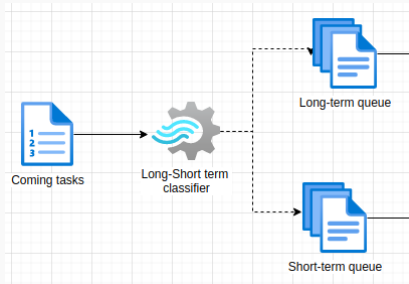


Figure: Long-short term classifier



Figure: Long term task scheduling

We should schedule long-term tasks first because in reality, they are always belonged to sensitive jobs such as Gmail, Docs in google trace data

## Scheduling algorithm

- ▶ Best fit: Optimize resources utilization
- ▶ Worst fit: Optimize time execution and workload balancing

# Proposed task scheduling pipeline

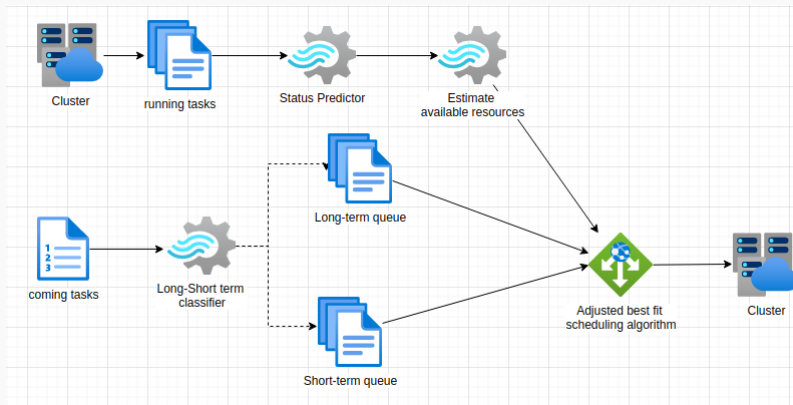


Figure: Task scheduling pipeline



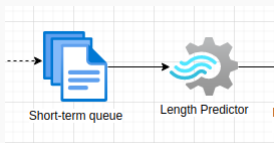


Figure: Length predictor

## Task information

- ▶ resources\_request
- ▶ meta\_data (priority, programming language, ...)
- ▶ length: the number of instructions

## Length estimation model

We expect that there is a machine learning model fitting task length data



## Task execution time fomula

$$T = \frac{\text{length}}{\text{MIPS} * \text{cpu\_usage}}$$

## Make-span

$$\text{makespan} = \max(T_1, T_2, \dots, T_n)$$

## Algorithm

**Input:** tasks, machines

**Output:** task-machine mapping

**Objective:** find a solution that minimize make-span of tasks



## Algorithm

**Algorithm 1** pseudocode for the calculation of

**Input:**  $L\_queue, S\_queue, vms$

**Output:** Map<Task, Vm>

- 1:  $usage \leftarrow \text{estimate\_usage}(vms)$
- 2:  $L\_running\_usage \leftarrow \text{get\_L\_running}(usage)$
- 3:  $S\_running\_usage \leftarrow \text{get\_S\_running}(usage)$
- 4:  $L\_solution \leftarrow \text{bestfit}(L\_queue, L\_running\_usage)$
- 5:  $L\_usage \leftarrow \text{update}(L\_running\_usage, L\_solution)$
- 6:  $S\_solution \leftarrow \text{bestfit}(S\_queue, S\_running\_usage)$
- 7: **return**  $\text{merge}(L\_solution, S\_solution)$

# Proposed task scheduling pipeline

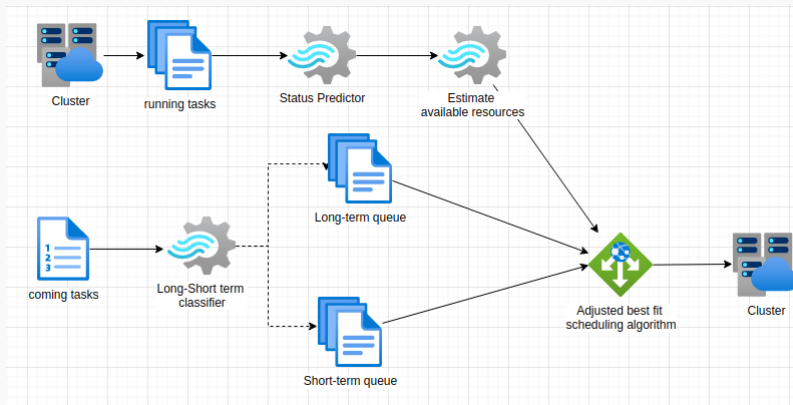


Figure: Task scheduling pipeline

# Update long-short term

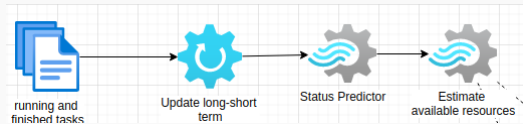


Figure: Update long-short term

When the duration of a short-term task is extremely longer than the expected, the task would be re-labeled as long-term

$$T = \frac{\text{length}}{\text{MIPS} * \text{cpu\_usage}} + \epsilon$$

# Update long-short term



## Task execution time

$$T \sim \mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{(T-\mu)^2}{2\sigma^2}}$$

## Task's duration

$$\begin{aligned} P(T > \text{duration}) &= \int_{\text{duration}}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{(t-\mu)^2}{2\sigma^2}} dt \\ &= \frac{1}{2} - \phi\left(\frac{\text{duration} - \mu}{\sigma}\right) \end{aligned}$$

## $\alpha$ threshold

$P(T > \text{duration}) < \alpha \rightarrow$  task is re-labeled as long-term

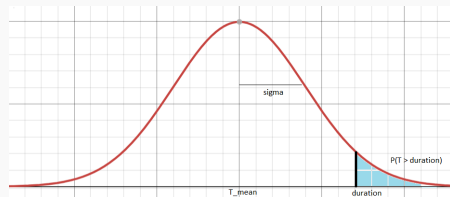


Figure: Task execution time distribution

## Uncertain state

- ▶ tasks are scheduled over  $state_1$
  - ▶ tasks are executed over  $state_2$
- $state_1$  and  $state_2$  are probably biased.

## Estimate $state_2$

Basing on  $state_1$ , we want to estimate  $state_2$  in order to get more accuracy environment information

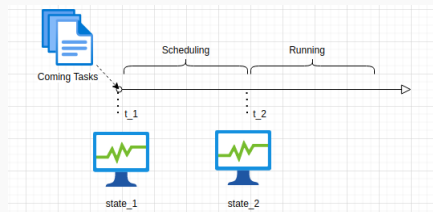


Figure: Scheduling process

## Task state

**Idea:** use Markov chain to estimate transition probabilities between the states

## Markov chain drawbacks

In experiment, the idea poses many problems :

- ▶ transition probabilities also depend on the duration of task
- ▶ cpu usage is also correlated to probability

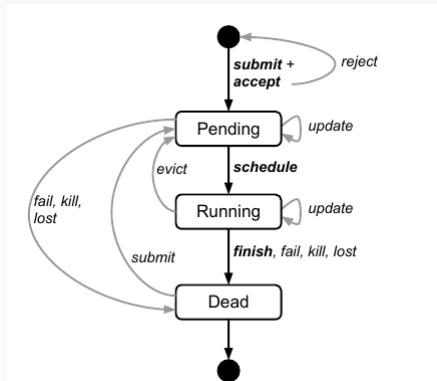


Figure: Task states transition





Figure: cpuRequest vs transition probability

Call:

```
lm(formula = prob ~ cpuRequest, data = df)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.304739	-0.124903	0.003684	0.117100	0.283832

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.12018	0.01649	7.287	1.59e-12 ***
cpuRequest	0.91220	0.08247	11.061	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1433 on 419 degrees of freedom  
Multiple R-squared: 0.226, Adjusted R-squared: 0.2242  
F-statistic: 122.3 on 1 and 419 DF, p-value: < 2.2e-16

Figure: Summary statistics

## Bayesian network

$$P(S_{i+1}|T_{i+1}, cpu, S_i) = \mathcal{F}_{S_{i+1}}(S_i, T_{i+1}, cpu, \theta)$$

## Probability function

$$\mathcal{A}_{S_{i+1}}(S, T, cpu, \theta) = \theta_0^{S_{i+1}} + \theta_1^{S_{i+1}} * S + \theta_2^{S_{i+1}} * T + \theta_3^{S_{i+1}} * cpu$$

$$\mathcal{F}_{S_{i+1}}(S_i, T_{i+1}, cpu, \theta) = \frac{\mathcal{A}_{S_{i+1}}(S, T, cpu, \theta)}{\sum_{S_k \in \{S\}} \mathcal{A}_{S_k}(S, T, cpu, \theta)}$$

state\state	pending	running	dead
pending	p11(duration, cpu)	p12(duration, cpu)	1 - p11 - p12
running	p21(duration, cpu)	p22(duration, cpu)	1 - p21 - p22
dead	p31(duration, cpu)	0	1 - p31

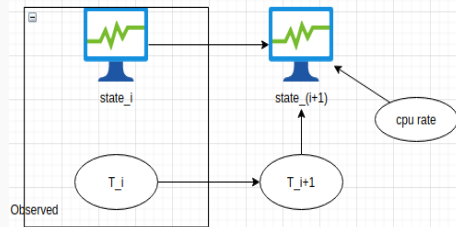


Figure: Transition network

## Probability factorization

$$\begin{aligned}P(S_{i+1}|S_i, T_i, cpu) &= \int P(S_{i+1}, T_{i+1}|S_i, T_i, cpu) dT_{i+1} \\&= \int P(S_{i+1}|S_i, T_{i+1}, cpu) P(T_{i+1}|T_i) dT_{i+1} \\&= E_{T_{i+1} \sim \mathcal{N}(T_i + \mu, \sigma^2)} \mathcal{F}_{S_i, cpu, \theta}(T_{i+1}) \\&\sim \mathcal{F}_{S_i, cpu, \theta}(T_i + \mu)\end{aligned}$$

Use gradient ascend to find  $\theta$  that maximize likelihood

$$\mathcal{L}(\theta : d) = \log(\mathcal{F}_{S_i, cpu, \theta}(T_i + \mu))$$

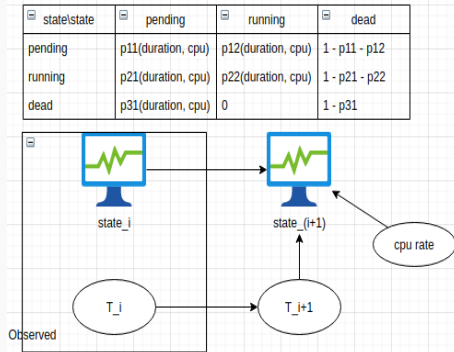


Figure: Transition network

# Resources estimation



## Tasks assigned to machines

$\mathcal{S}$  is a set of tasks which are assigned to the machine and haven't finished

## Resources usage estimation

$$resource\_usage = \sum_{task \in \mathcal{S}} P(running) * resource\_usage\_of\_task$$

$$available\_resource = total\_capacity - resource\_usage$$

- $P(running)$ : probability that the task continue running after scheduling time

## Scheduling information

Coming tasks would be scheduled over resources estimation, not observed resources at scheduling time

# Table of Contents



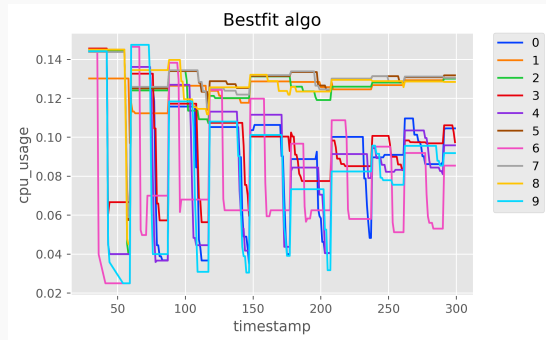
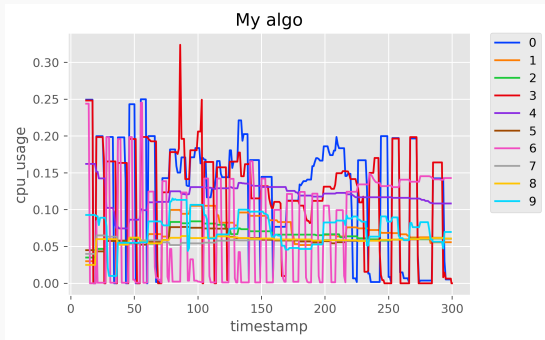
System overview

Problem Statement

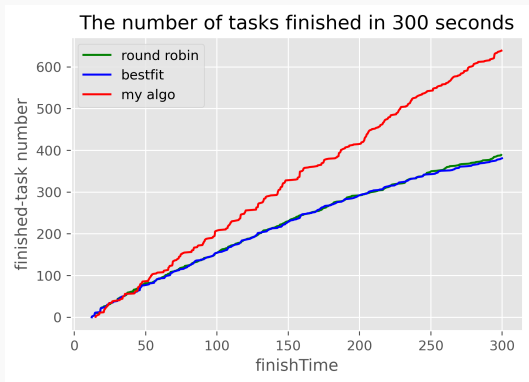
Solution

Experiment result

# Cpu usage experiment

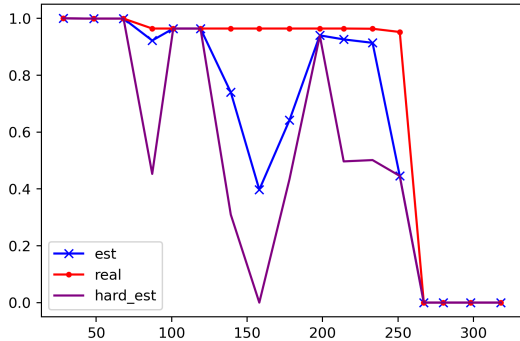


# Number of tasks finished

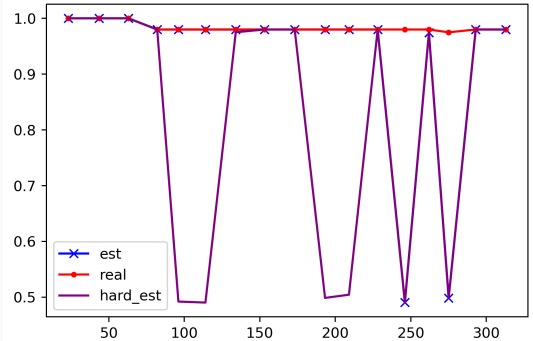


Statistics	RoundRobin	BestFit	MyAlgo
count	390.00	382.00	640.00
mean	28.14	24.12	13.15
std	26.74	22.77	15.51
min	1.99	2.00	2.62
25%	8.83	7.83	5.38
50%	21.25	18.07	7.23
75%	36.46	31.19	14.82
max	148.15	131.54	109.16

vm6 cpu usage estimation



vm9 cpu usage estimation

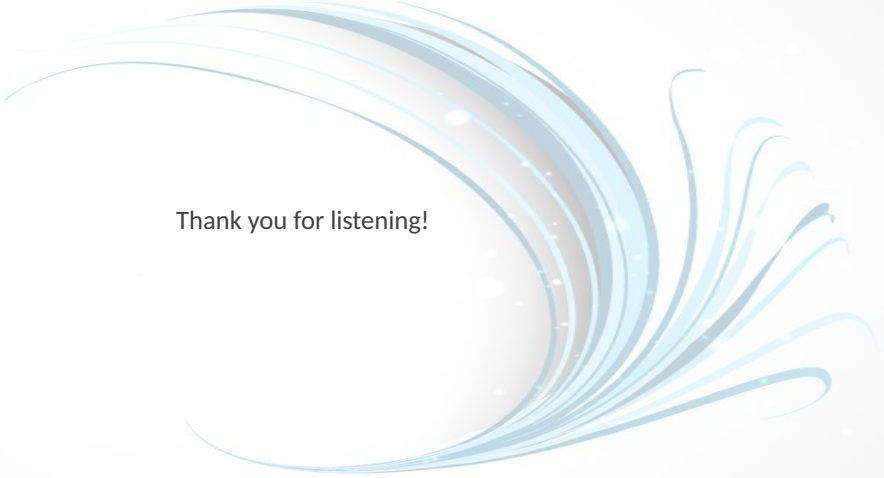






## Comparison between available-resources estimation vs no estimation

	<b>statistic</b>	<b>no-estimated</b>	<b>estimated</b>
<b>0</b>	count	366.000000	378.000000
<b>1</b>	mean	32.177039	30.375850
<b>2</b>	std	23.648688	23.509136
<b>3</b>	min	6.644639	6.581494
<b>4</b>	25%	17.068980	15.700250
<b>5</b>	50%	27.211863	24.485500
<b>6</b>	75%	37.649000	33.975500
<b>7</b>	max	141.561000	138.182000



Thank you for listening!