

Energy efficiency in multi-core supercomputing

Tuan V. Dinh

Supervisors: A/Prof. Lachlan Andrew, Dr. Philip Branch and Dr. Yoni Nazarathy



Motivations



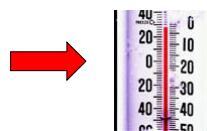
Energy consumption in supercomputer



(1)



Electricity bills



Heat – thermal management



Investment – cooling systems, hardware, etc.

Motivations



- Trend: desired for more powerful computing,
 - climate change, astronomy
 - “exascale” supercomputing
- Graphics Processing Units to replace ordinary CPUs



(1) Green Machine



(2) GStar – Green II

SWINBURNE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

(1),(2) <http://astronomy.swin.edu.au/supercomputing/>

CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 3

Prior work



Design and thermal management.

(Al-Fares et. al 2008, Sharma et. al, 2005; Patel et. al, 2002...)

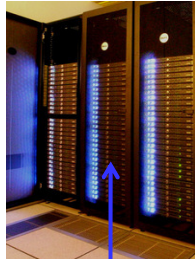
SWINBURNE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 4

Prior work



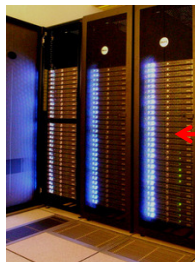
**Energy efficiency in
single processor**
(Yao 1995; Wierman et. al 2009;...)



CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 5

Prior work



**Workload scheduling,
resource consolidation.**
(Lin et al. 2007; Bundle
2006;Pinherio et. al. 2001;...)



CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 6

Prior work



**Capacity Provisioning
(Power Proportionality)**
(Lin et. al 2011; Xiong 2010; Chase et al. 2001;...)



**Workload scheduling,
resource consolidation.**
(Lin et al. 2007; Bundle 2006; Pinheiro et. al. 2001;...)



CAIA Seminar

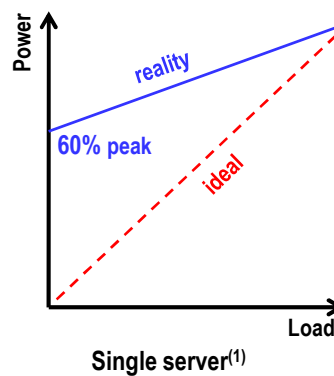
<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 7

Power proportionality



**Capacity Provisioning
(Power Proportionality)**

- When idle, server still consumes at least 60% of peak power
- Save large amount of energy by powered them off
- Challenges:
 1. Short of supply in near future
 2. Increase hardware failure rate



(1) Bassoro, "The case for energy proportional", 2007.

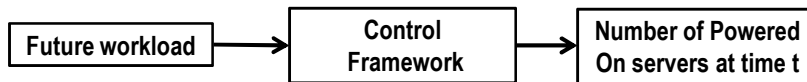
CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 8

Research question



Capacity Provisioning
(Power Proportionality)



SWIN
BUR
NE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 9

Outline



- Control framework
- Future load estimation
- Energy efficiency for single core
- Future work and challenges

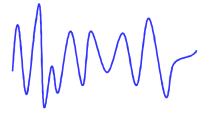
SWIN
BUR
NE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 10

Control framework



Future arrival load $A(t)$:

Number of active server at time t

$$\min_{n(0), n(1), \dots} \sum_{t=0}^{\infty} c o s$$

Server "On" > A

Cost = Energy (E) + Switching cost (SC) + Performance degradation cost (PC)

\propto number of servers

\propto number of switching servers

\propto times of short of supply

Challenges:

- Long future load estimation horizon
- Intractable for implementation

SWINBURNE

SWINBURNE UNIVERSITY OF TECHNOLOGY

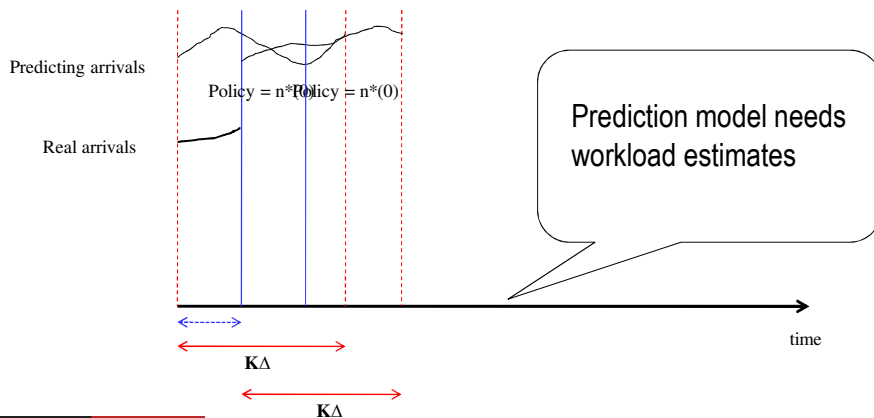
CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 11

Model Predictive Control (MPC)



- MPC is a control method for stochastic control problem



SWINBURNE

SWINBURNE UNIVERSITY OF TECHNOLOGY

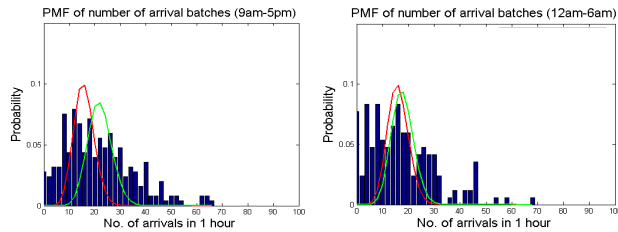
CAIA Seminar

<http://caia.swin.edu.au/cv/tdinh> 22 August 2012 Slide 12

Green Machine historical workload



- Our experiments show that Green Machine's workload is not Poisson



- Characteristics:

- Jobs arrive in "batch". (with similar runtime and cluster size)
- The arrival intensity rate is varying

$M_x^t / G / \infty$ with Model Predictive Control



- $A(t)$ is modelled as a batch Poisson process with time varying rate (M_x^t)
- Service time distribution is stationary
- Sufficient supply

Only consider
K slots:

$$\min_{n(0), n(1), n(K-1), \dots} \sum_{k=0}^{K-1} g(k)$$

- Only use $n^*(0)$ and discard the rest
- Use the same process at the next slot, i.e for slot $k=1$ to $k = K+1$

Case A



Cost(t) = "On" Servers + server switching + Not meeting the demand

$$\text{cost}(t) = n(t) + u(t) + \mathbb{E}[(X(t) + U(t) - n(t))^+]$$

jobs that arrived after t = 0 and still running

jobs that arrived before t = 0 and still running

PD cost

Case B



Cost(t) = "On" Servers + server switching

$$\text{cost}(t) = n(t) + u(t)$$

No. "ON" servers

No. "switching" servers

but apply the constraint:

$$P\{X(t) + U(t) \geq n(t)\} < 1 - \epsilon$$

PD cost is decoupled to be a constraint

M_x^t /G/∞ with Model Predictive Control



Case A: $\min_{n(0), n(1), n(k-1), \dots} \sum_{k=0}^{K-1} g(k)$	Case B: $\min_{n(0), n(1), n(k-1), \dots} \sum_{k=0}^{K-1} \beta_1 n(k) + \beta_2 u(k)$ s. t $\Pr[\{x(t) + u(t)\} < n(k)] > 1 - \epsilon$
<ul style="list-style-type: none"> • Unconstrained optimisation problem • PD cost is integrated in objective cost • Can be cast as shortest path problem and solved by Dynamic Programming 	<ul style="list-style-type: none"> • Constrained optimisation problem • Solved using Linear Programming • The performance constraint is decoupled from objective cost • Cost is deterministic • Give explicit bound of performance (ϵ)

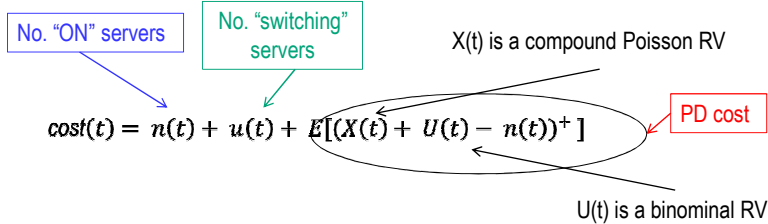


SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 17

Solving case A:



- Use Gaussian approximations $\bar{X}(t)$ and $\bar{U}(t)$
- Cost(t) can then be calculated based on:
 - service time distribution, batch size distribution
- As cost per stage cost(t) is now determined, use Dynamic Programming



SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 18

Outline



- Control framework

- Future load estimation

- Energy efficiency for single core

- Future work and challenges

Historical workload study



- Concentrate on the arrival process $A(t)$
 - estimating the future load
- Swinburne Supercomputer maintains accounting logs
 - keep job records (arrival time, runtime, cluster size, used resources, etc.)
 - guidance for workload modelling, simulation and validation.
- Attempts:
 - statistical models for interarrival time, batch size, runtime and arrival intensity
 - clustering users based on job statistics

Future workload estimation

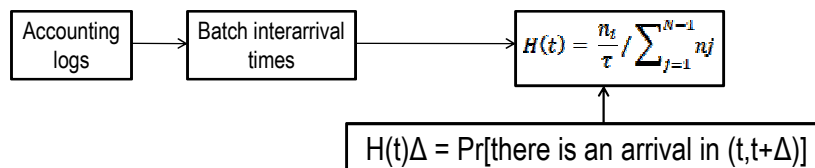


- Users have different submission behaviours
 - but it seems consistent for each user.
- Is knowing per user information helpful ?
 - method to predict future load
 - hazard rate function of interarrival time distribution
 - how to testify that ?
 - real trace experiment

Future workload estimation



- Hazard rate function of the interarrival time distribution:

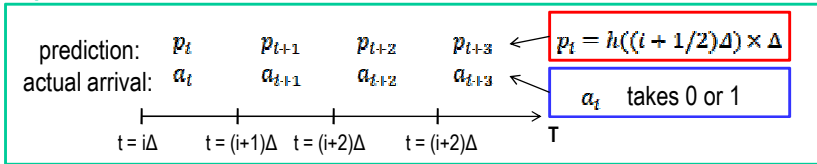


- The likelihood of having an arrival in short future, provided the time of the last submission

Evaluating the load estimator



Experiment:



Length/Scenario	Not knowing user behaviours $\frac{\Delta}{T} \sum_i^{T/\Delta} (p_i^1 - a_i)^2$	Knowing per user behaviours $\frac{\Delta}{T} \sum_i^{T/\Delta} (p_i^2 - a_i)^2$
6 months	0.2207	0.1559
10 months	0.2114	0.1993

- ~ 5% improvement when knowing user behaviours
- the fitting is not perfect
- only predict the arrival **likelihood**
- the ultimate prediction is number of servers needs to be powered ON

- Next steps:**
- Predict the required servers
 - How?: required servers = likelihood x average batch size



CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 23

Outline



- Control framework
- Future load estimation
- Energy efficiency for single core
- Future work and challenges



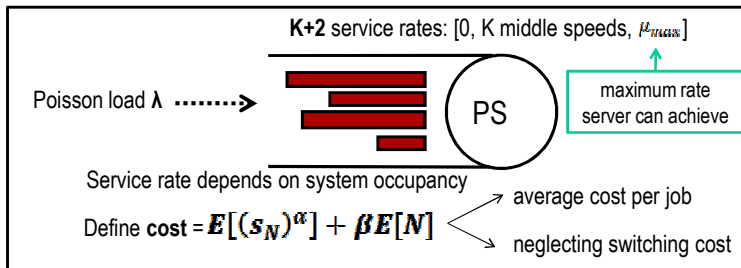
CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 24

Energy efficiency in single core



M/G/1-PS:



Research questions:

1. Knowing system load λ_d , which service rates to use ?
2. The benefit of having more choices of rates (larger K)

SWINBURNE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 25

Energy efficiency in single core



Prior work	Our contributions
<ol style="list-style-type: none"> 1. Not much cost improvement when $K = \infty$ comparing when server only runs at one optimal rate, if $\lambda_a = \lambda_d$. [Wierman. et al. 2009] 2. When $\lambda_a \neq \lambda_d$, server with $K = \infty$ is more robust than server has only one (optimal) rate, in terms of cost variation. [Andrew. et al., 2010] 	<ol style="list-style-type: none"> 1. Given a λ_d, find the optimal rates for server with $1 \leq K < \infty$ 2. Our numerical results suggest that the local robustness (small variation) increases as K increases. 3. $K = 8$ is as good as $K = \infty$ in terms robustness 4. Flexibility in choosing rate in runtime gives significant cost improvement

SWINBURNE

SWINBURNE
UNIVERSITY OF
TECHNOLOGY

CAIA Seminar

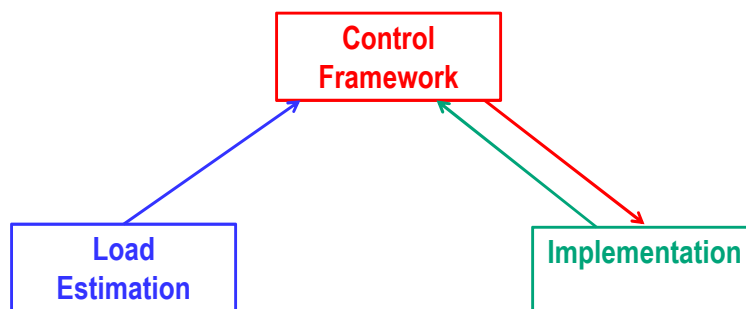
<http://caia.swin.edu.au/cv/dinh> 22 August 2012 Slide 26

Outline



- Control framework
- Future load estimation
- Energy efficiency for single core
- Future work and challenges

Future work



Control Framework – Future work



Control Framework

SOLVING ALL OPTIMIZATIONS

- Penalty of using Gaussian approximations
- Find suitable model for service time distribution (apply for both cases)

SIMULATION

- An event driven simulator has been written
 - Using the Green Machine scheduling policy
- Compare performance of the controller to the “halfway-house” heuristic
 - If found no improvement, examining where the energy savings is claimed.

Future work



Load Estimation

- Instead of predicting arrival likelihood, predict the requesting number of servers
 - still use the hazard rate function
 - mean batch size for each user
- Modelling the service time distribution (find the c.d.f) and the batch size distribution
 - starting with Weibull, Gamma and Log-Normal distribution

Integrating the controller



- Challenges:
 - Not all machines can be powered OFF: storage, special reservation
 - The scheduler Moab is a commercial software
 - Require to understand Moab decision
 - But Moab logs are accessible
- Future tasks (time permitting)
 - Integrate the load estimation into the controller
 - Develop the mechanism to turn server ON and OFF
 - Not interfere with normal operations
 - Develop interfaces between the controller and Moab

Implementation

Summary



- Power proportionality can achieve much energy saving for supercomputer.
- A control framework has been developed using Model Predictive Control, but it needs to be tested.
 - A controller is developed based on the framework
- An estimation of future workload is required in the control framework
 - Knowing per user behaviour can improve the estimation accuracy