

On Simulation and Design of Parallel-Systems Schedulers

Are We Doing the Right Thing?

Xinyu Chen

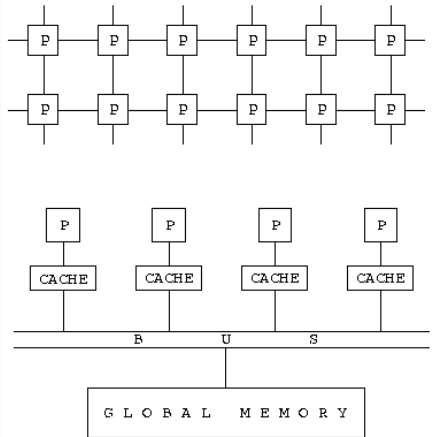
April 13, 2017

University of New Mexico

1. A Little Background
2. Scheduling Simulator
3. Experiments Result
4. Conclusion and Related Work
5. Questions

A Little Background

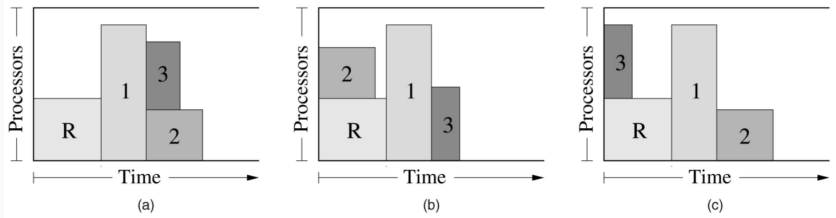
PARALLEL SYSTEM ARCHITECTURE



- Distributed-memory model
- shared-memory model

Scheduling Simulator

CONVENTIONAL SCHEDULERS AND SIMULATIONS



Trace-driven Workload Some Metrics

- Close-System
- Open-System
- Response Time = $t_{\text{terminate}} - t_{\text{start}}$
- Slowdown = $\frac{\text{resp_time}}{\text{actual_runtime}}$
- Thinking Time = $t_{\text{new_submission}} - t_{\text{last_finish}}$

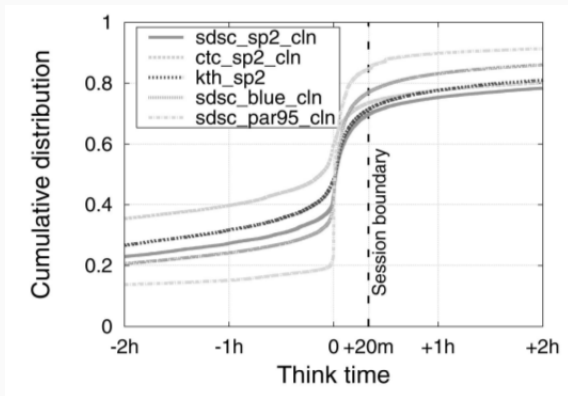
Table: Methodological Difference between Two Types of Simulation

Category	Conventional Simulations	Site-Level Simulations
Workload source	System traces	User models
Workload generation	Open-system model	User-scheduler interaction
Load scaling	Trace (de)-compression	Number of users
Performance Metrics	Response time, slowdown	Throughput, session length

User Model

- Session Dynamic
- Job Submission
- Activity Cycle

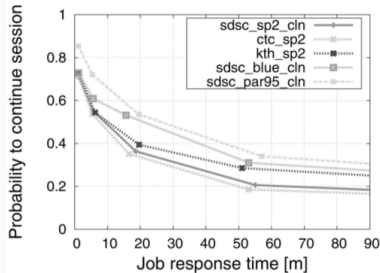
USER BEHAVIOR - FROM THINK TIME TO SESSIONS



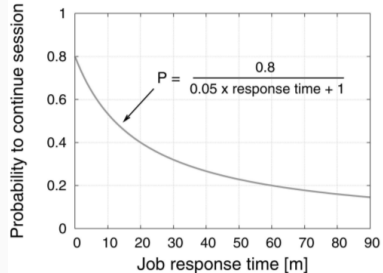
Potential Problems

The CDF of think time looks symmetric. The negative part indicates script submit.

USER BEHAVIOR - SESSION DYNAMICS



(a)



(b)

Probability to continue session

$$p_{\text{cont}}(j) = \frac{0.8}{0.05 \times \text{resp_time}(j) + 1}$$

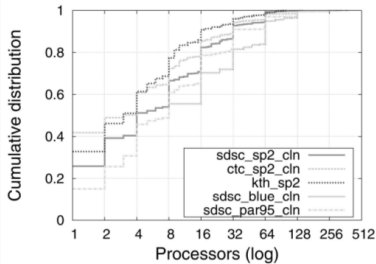
Potential Problems

Regression?

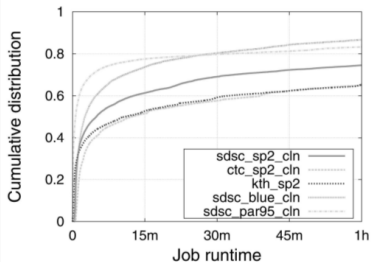
Confounding Factors?

Causal Relationship?

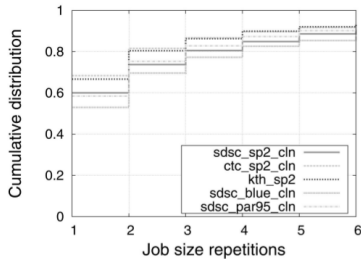
USER BEHAVIOR - JOB SUBMISSION



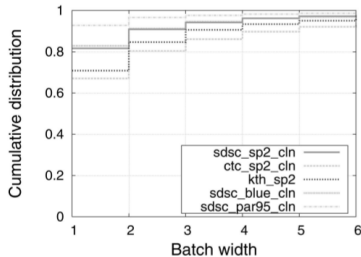
(a)



(b)



(a)



(b)

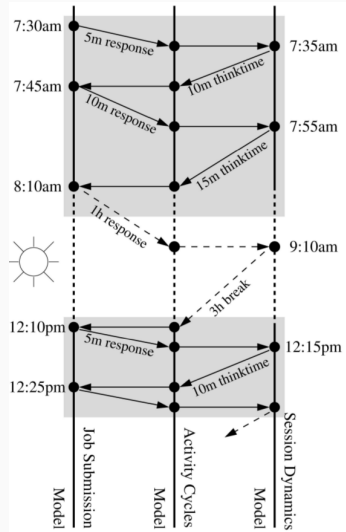
USER BEHAVIOR - ACTIVITY CYCLE

Table: Four user cycle classes

Daytime-weekdays	Daytime-weekend
Nighttime-weekdays	Nighttime-weekend

Recap User Model

- Workload source
- Workload Generation



Similar to EASY

$$criticality(j) = \frac{0.04}{(0.05 \times estimate_resp_time(j) + 1)^2}$$

$$priority(j) = \alpha \times criticality(j) + seniority(j)$$

$$estimate_resp_time(j) = seniority(j) + runtime(j)$$

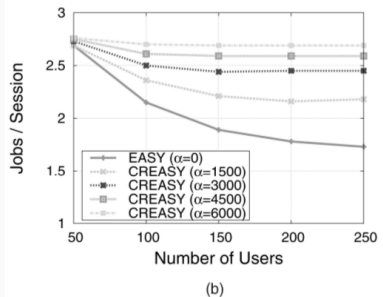
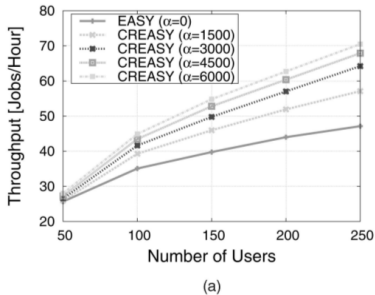
Potential Problems

criticality tends to starve long jobs

How to decide and interpret α values

Experiments Result

METRICS: THROUGHPUT AND SESSION LENGTH



Potential Problems

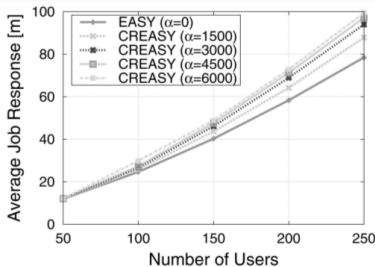
Is Number of Users related with Productivity?

How many processors are used?

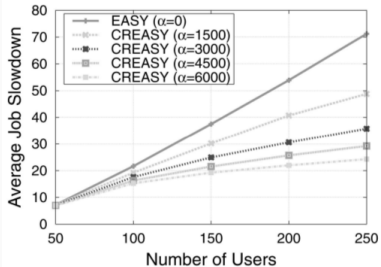
What are the job sizes?

Larger number of users favors large α , which favors short jobs.

METRICS: RESPONSE TIME AND SLOWDOWN



(a)



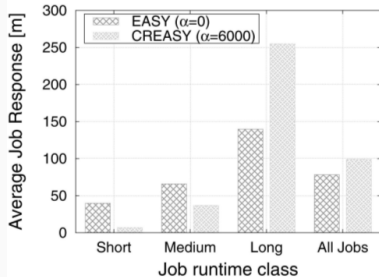
(b)

Potential Problems

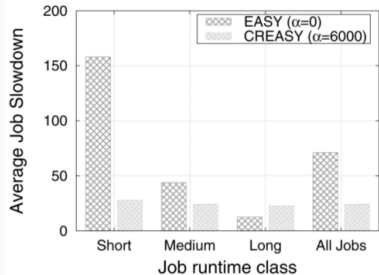
Is Average Response Time and Slowdown useful? What is the standard error?

Slowdown should be greater than 100%.

METRICS: RESPONSE TIME AND SLOWDOWN IN DIFFERENT CLASSES



(a)



(b)

Potential Problems

It's clear in the above two graphs that the scheduler favors small jobs and sacrifices big jobs.

Slowdown should be greater than 100%

Conclusion and Related Work

CONCLUSION

- Trace-driven Scheduling Simulation cannot reflect user-system interactions.
- Conventional metrics like Response Time and Slowdown cannot help improve user productivity
- The User Model is much simplified. Real user behaviors are more complicated.
- The interactive sessions is critical but only require less resources. Script submission need to be considered in the future.

SOME PREVIOUS DISCUSSED PAPERS

- Looking at Data by Dror G. Feitelson
- How Uber Uses Psychological Tricks to Push Its Drivers' Buttons

Questions
