

COMP9334

# Capacity Planning for Computer Systems and Networks

---

Week 9b: Further applications of queueing

# Applications of queueing

---

- There are plenty and we will look at a few examples
  - The technical papers can be downloaded from the course website (password required)
- Good resource:
  - Sigmetrics
    - <http://www.sigmetrics.org>
    - A leading conference on performance evaluation of computer systems and networks
  - The journal Performance Evaluation

# Determining Multi-programming level

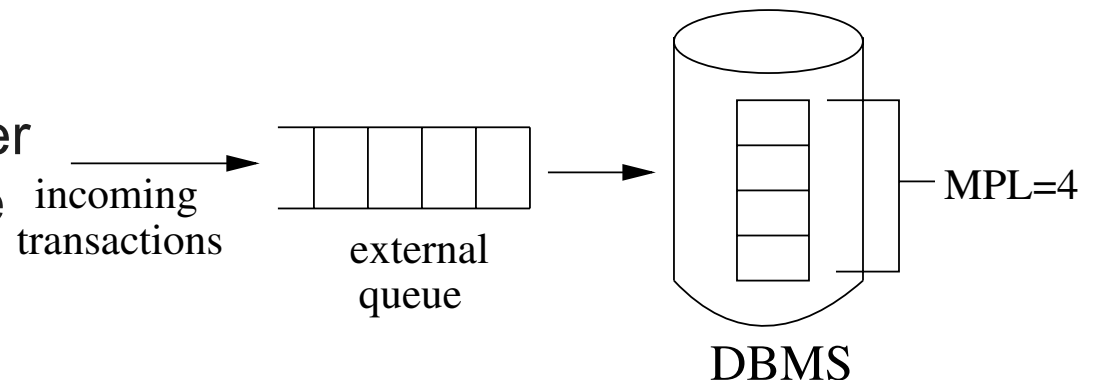
---

## How to determine a good multi-programming level for external scheduling

Bianca Schroeder <sup>§</sup>	Mor Harchol-Balter <sup>§*</sup>	Arun Iyengar <sup>†</sup>	Erich Nahum <sup>†</sup>	Adam Wierman <sup>§</sup>
§Carnegie Mellon University		†IBM T.J. Watson Research Center		
Department of Computer Science		Yorktown Heights, NY USA		
Pittsburgh, PA USA		<aruni,nahum>@us.ibm.com		
<bianca, harchol, acw>@cs.cmu.edu				

# DB server – Multi-programming level

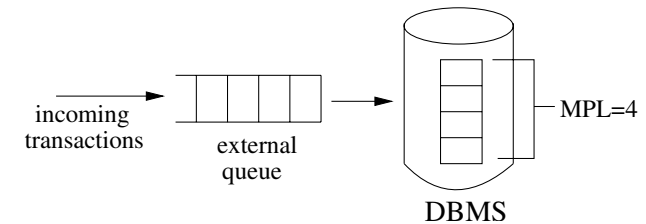
- Some database server management systems (DBMS) set an upper limit on the number of active transactions within the system
- This upper limit is called multi-programming level (MPL)



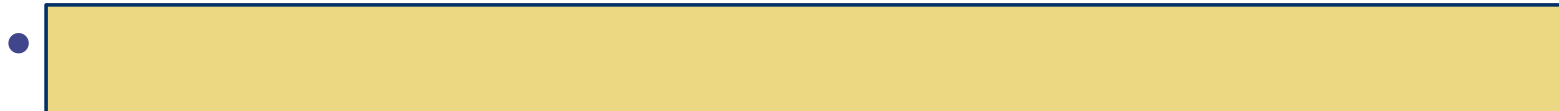
**Figure 1.** *Simplified view of the mechanism used in external scheduling. A fixed limited number of transactions ( $MPL=4$ ) are allowed into the DBMS simultaneously. The remaining transactions are held back in an external queue. Response time is the time from when a transaction arrives until it completes, including time spent queueing externally to the DBMS.*

- A help page from SAP explaining MPL
- [http://dcx.sap.com/1200/en/dbadmin\\_en12/running-s-3713576.html](http://dcx.sap.com/1200/en/dbadmin_en12/running-s-3713576.html)
- Picture from Schroder et al. “How to determine a good multi-programming level for external scheduling”

# The problem



- To choose a good MPL means you want to determine the mean response time for different choices of MPL
  - If  $MPL = 1$ , what is the response time?
  - If  $MPL = 2$ , what is the response time?
  - ...
- Question: Let us assume that the arrival is Poisson, can you suggest how we can determine the mean response time?



## Optimal Power Allocation in Server Farms

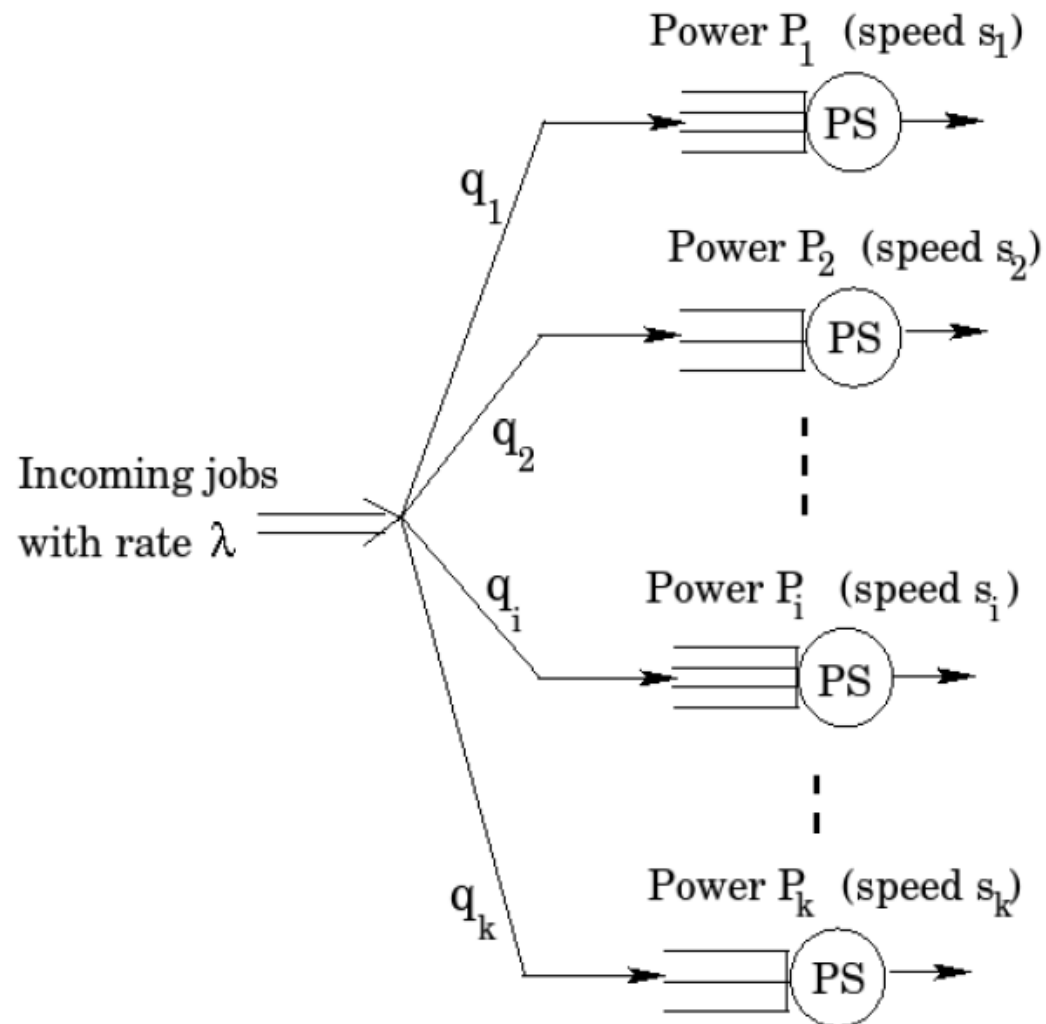
Anshul Gandhi  
Carnegie Mellon University  
Pittsburgh, PA, USA  
anshulg@cs.cmu.edu

Rajarshi Das  
IBM Research  
Hawthorne, NY, USA  
rajarshi@us.ibm.com

Mor Harchol-Balter\*  
Carnegie Mellon University  
Pittsburgh, PA, USA  
harchol@cs.cmu.edu

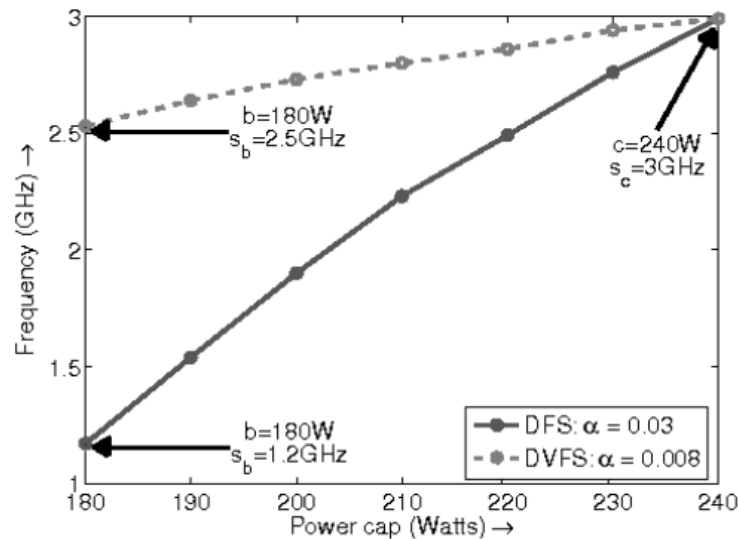
Charles Lefurgy  
IBM Research  
Austin, TX, USA  
lefurgy@us.ibm.com

# Server farm model

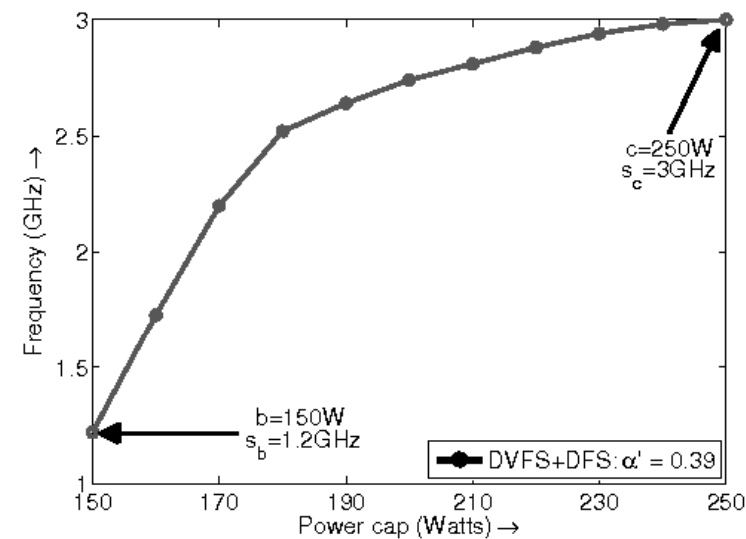


PS stands for  
processor sharing

# Power-frequency curve



(a) DFS and DVFS



(b) DVFS+DFS

*Power-to-frequency curves for DFS, DVFS, and DVFS+DFS for the CPU bound LINPACK workload. Fig.(a) illustrates our measurements for DFS and DVFS. In both these mechanisms, we see that the server frequency is linearly related to the power allocated to the server. Fig.(b) illustrates our measurements for DVFS+DFS, where the power-to-frequency curve is better approximated by a cubic relationship.*

PS stands for  
processor sharing



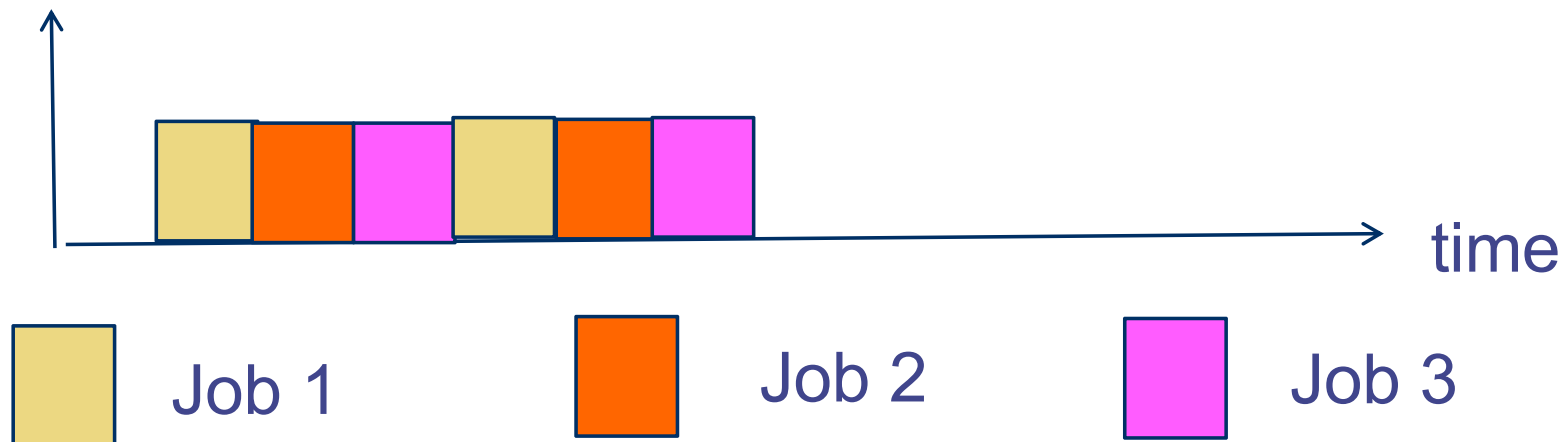
# The power allocation problem

---

- A server farm consists of multiple servers
- The servers can run at
  - Higher clock speed with higher power
  - Lower clock speed with lower power
- Ex: Given
  - Higher power = 250W, lower power = 150W
  - Power budget = 3000W
  - You can have
    - 12 servers at highest clock speed
    - 20 servers at lowest clock speed
    - Other combinations
  - Which combination is best?

# Processor sharing (PS)

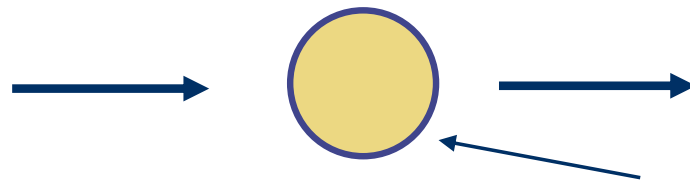
- For many operating systems, the processor works on a job for a quanta and then switch to another job for another quanta



## M/G/1/PS

- Poisson arrivals with mean arrival rate  $\lambda$
- General service time distribution with mean rate  $\mu$
- Processing sharing (PS)
- Mean response time

$$= \frac{1}{\mu - \lambda}$$



**Processor sharing**

# Power allocation for 2 servers

---

- To be worked out during the lecture

# Conclusions

---

- Queueing theory has many applications
- You have learnt the basics of analysis and simulation
- There are a lot of advanced theory and methods that we cannot cover but the basics will enable you to learn more