

COMP9334

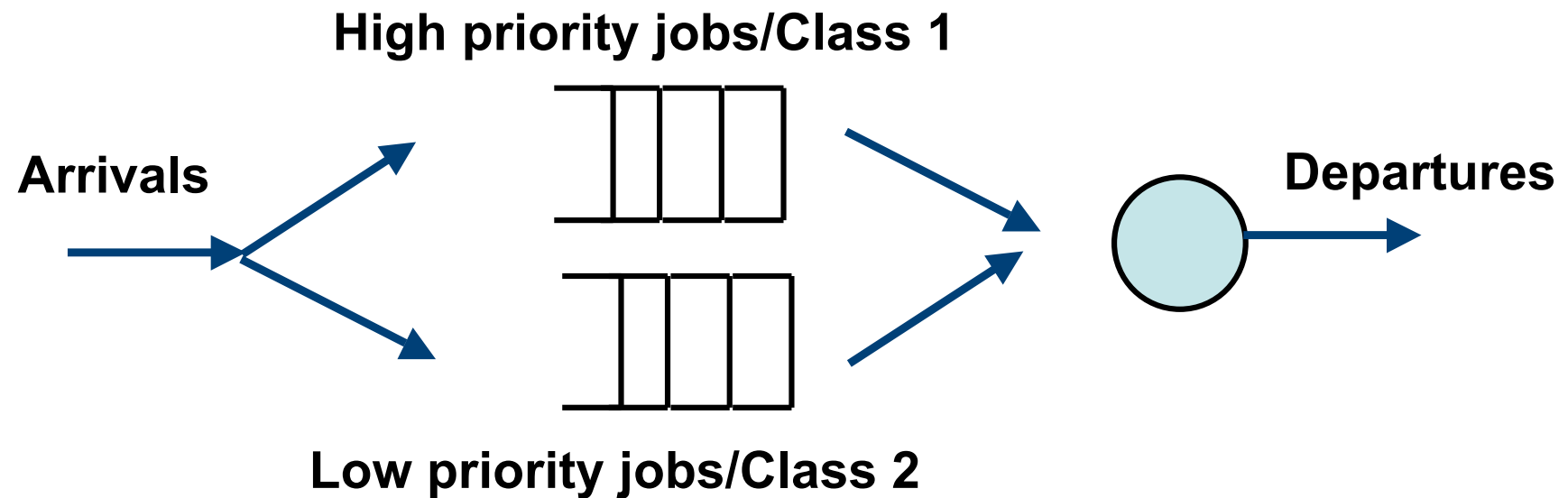
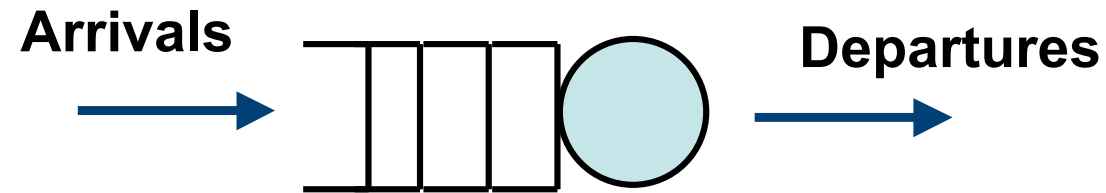
Capacity Planning for Computer Systems and Networks

Week 6B_2: Using queueing theory to improve system performance

Applications of queueing

- There are plenty of examples on using queueing to improve system performance
- We will look at a few examples
 - The technical papers can be downloaded from the course website
- Good resource:
 - Sigmetrics
 - <http://www.sigmetrics.org>
 - A leading conference on performance evaluation of computer systems and networks
 - Performance evaluation
 - A journal devoted to the topic of performance evaluation

Priority queueing



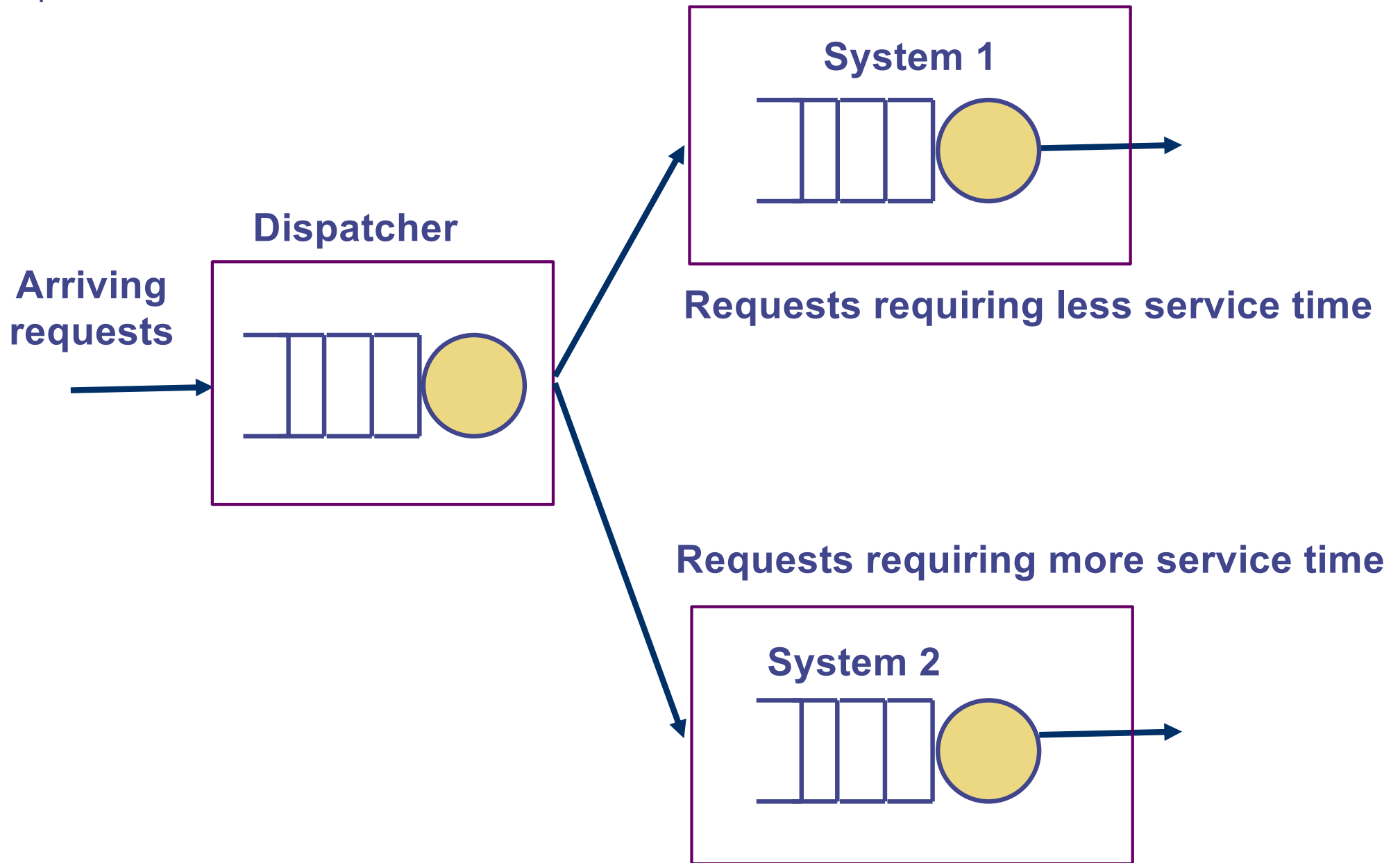
Using priority queueing to reduce response time (1)

- Revision Problem 4A, Question 1
 - Customers arrive at a grocery store's checkout counter according to Poisson process with rate 1 per minute.
 - Each customer carries a number of items that is uniformly distributed between 1 and 40.
 - The store has 2 checkout counters, each capable of processing items at a rate of 15 per minute.
 - To **reduce** customer waiting time in queue, the store manager considers dedicating one of the two counters to customers with x items or less and dedicating the other counter to customers with more than x items.
 - Write a small computer program to find the value of x that minimises the average customer waiting time.

Using priority queueing to reduce response time (2)

- At a website, use a dispatcher to classify the HTTP requests depending on their service time requirement
 - Requests whose service time is below a threshold go to one server
 - The rests of the requests go to the other server
 - See diagram on the next slide
 - Need prior knowledge on the service time of the requests
-
- Reference: B. Schroeder and M. Harchol-Balter, "Web servers under overload: How scheduling can help", ACM Transactions on Internet Technology (TOIT), Volume 6 , Issue 1, pages 20 - 52, 2006. <http://doi.acm.org/10.1145/1125274.1125276>.

Dispatcher



Determining Multi-programming level

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

Bianca Schroeder [§]	Mor Harchol-Balter ^{§*}	Arun Iyengar [†]	Erich Nahum [†]	Adam Wierman [§]
§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)

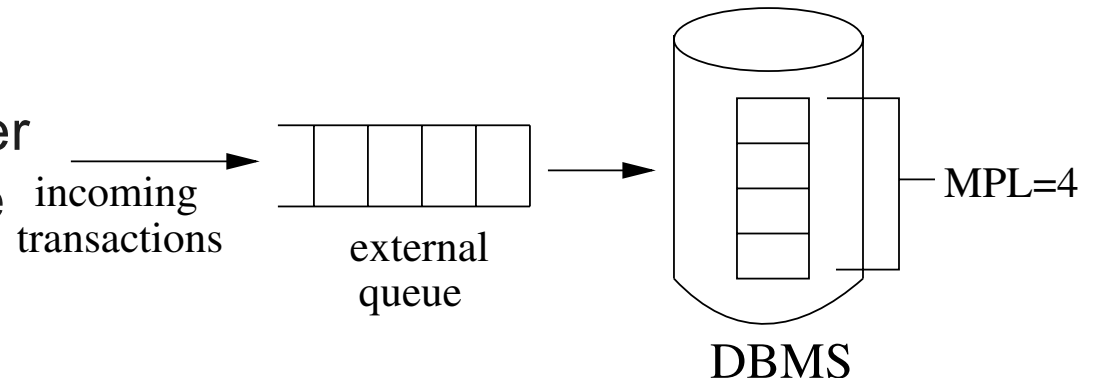
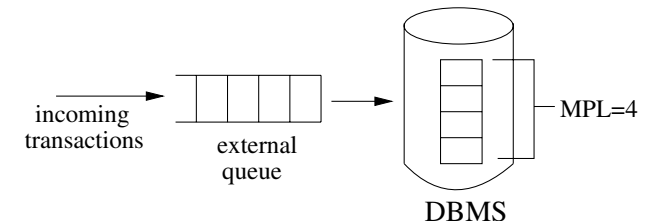



Figure 1. *Simplified view of the mechanism used in external scheduling. A fixed limited number of transactions (MPL=4) are allowed into the DBMS simultaneously. The remaining transactions are held back in an external queue. Response time is the time from when a transaction arrives until it completes,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

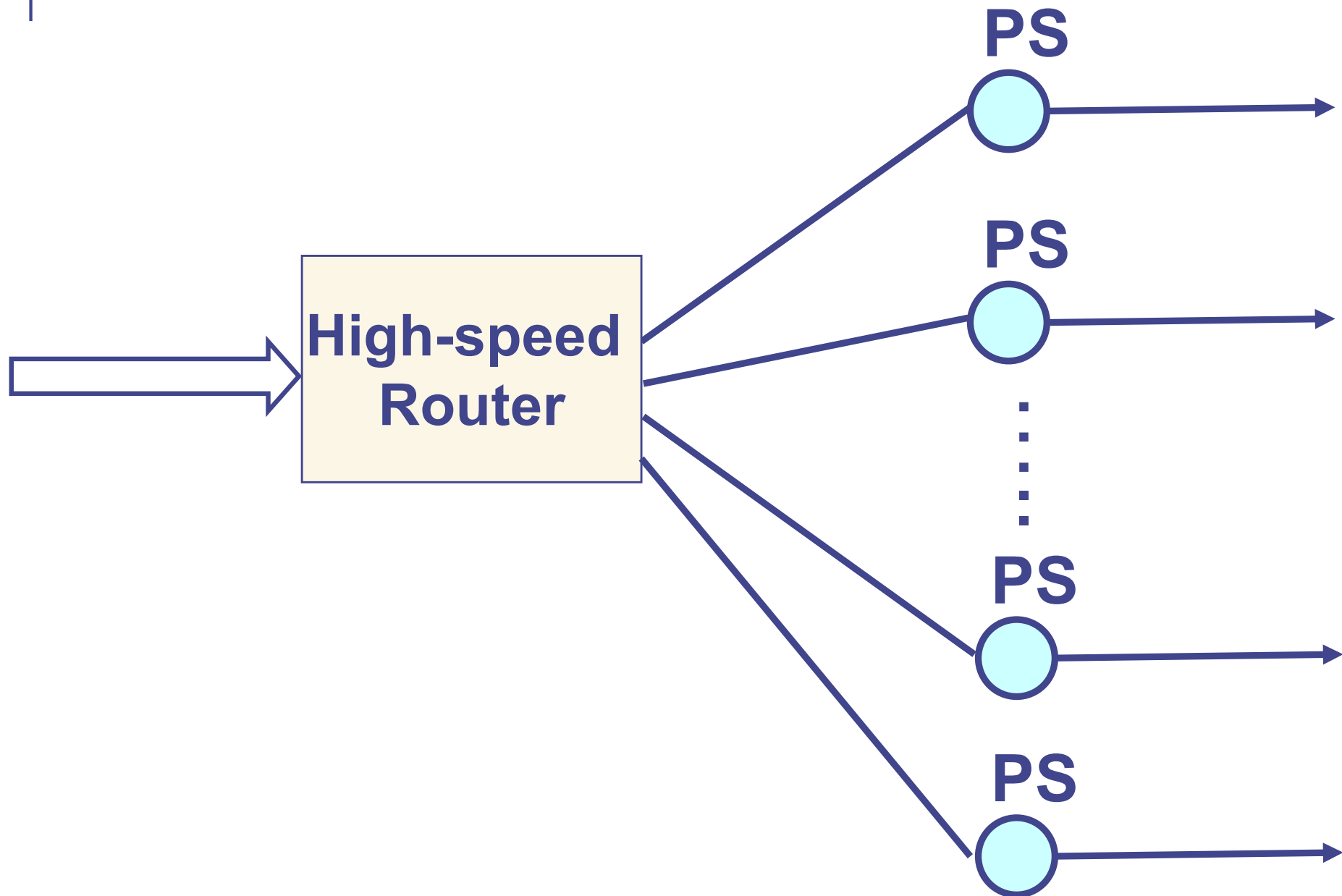
Rajarshi Das
IBM Research
Hawthorne, NY, USA
rajarshi@us.ibm.com

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

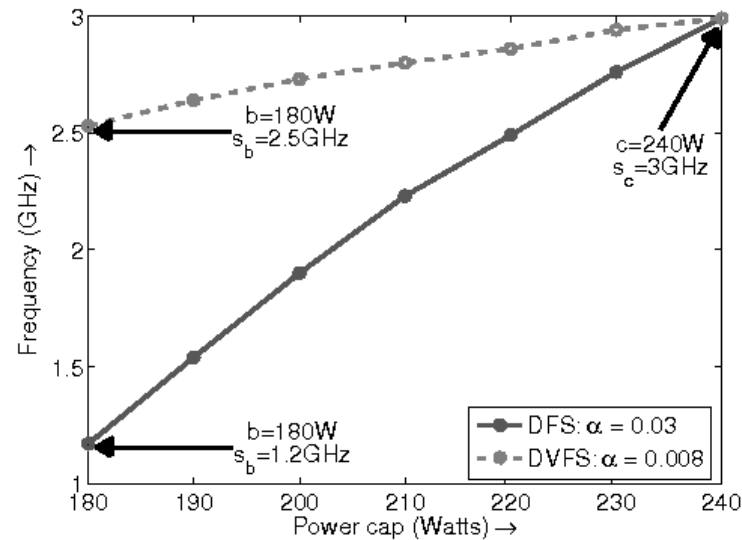
Charles Lefurgy
IBM Research
Austin, TX, USA
lefurgy@us.ibm.com

Server farm power allocation: Introduction

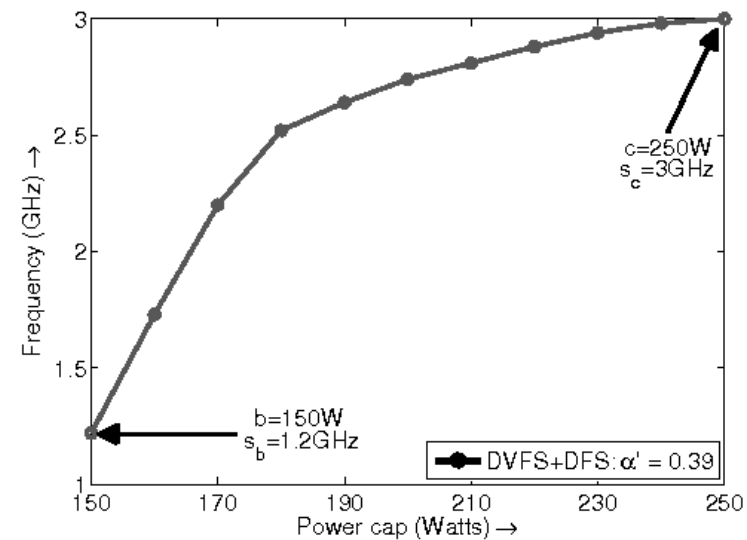
- A server farm consists of multiple servers
 - See next slide
- The servers can run at
 - Higher clock speed with higher power
 - Lower clock speed with lower power
 - Illustration: The slide after next
- Ex: Given
 - Higher power = 250W, lower power = 150W
 - Power budget = 3000W
 - You can have
 -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



CPU Power-speed tradeoff



(a) DFS and DVFS



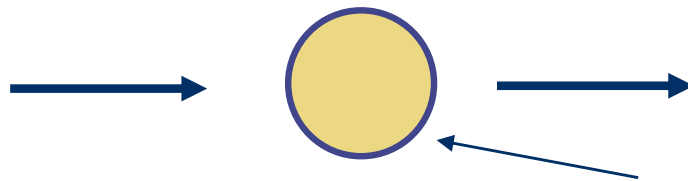
(b) DVFS+DFS

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.

M/G/1/PS

- Poisson arrivals with mean arrival rate λ
- General service time distribution with mean rate μ
- Processing sharing (PS)
- Mean response time

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



Processor sharing

Power allocation for 2 servers

- To be worked out during the lecture

Server farms with setup costs

Performance Evaluation 67 (2010) 1123–1138



Contents lists available at ScienceDirect

Performance Evaluation

journal homepage: www.elsevier.com/locate/peva



Server farms with setup costs

Anshul Gandhi^{a,*}, Mor Harchol-Balter^a, Ivo Adan^b

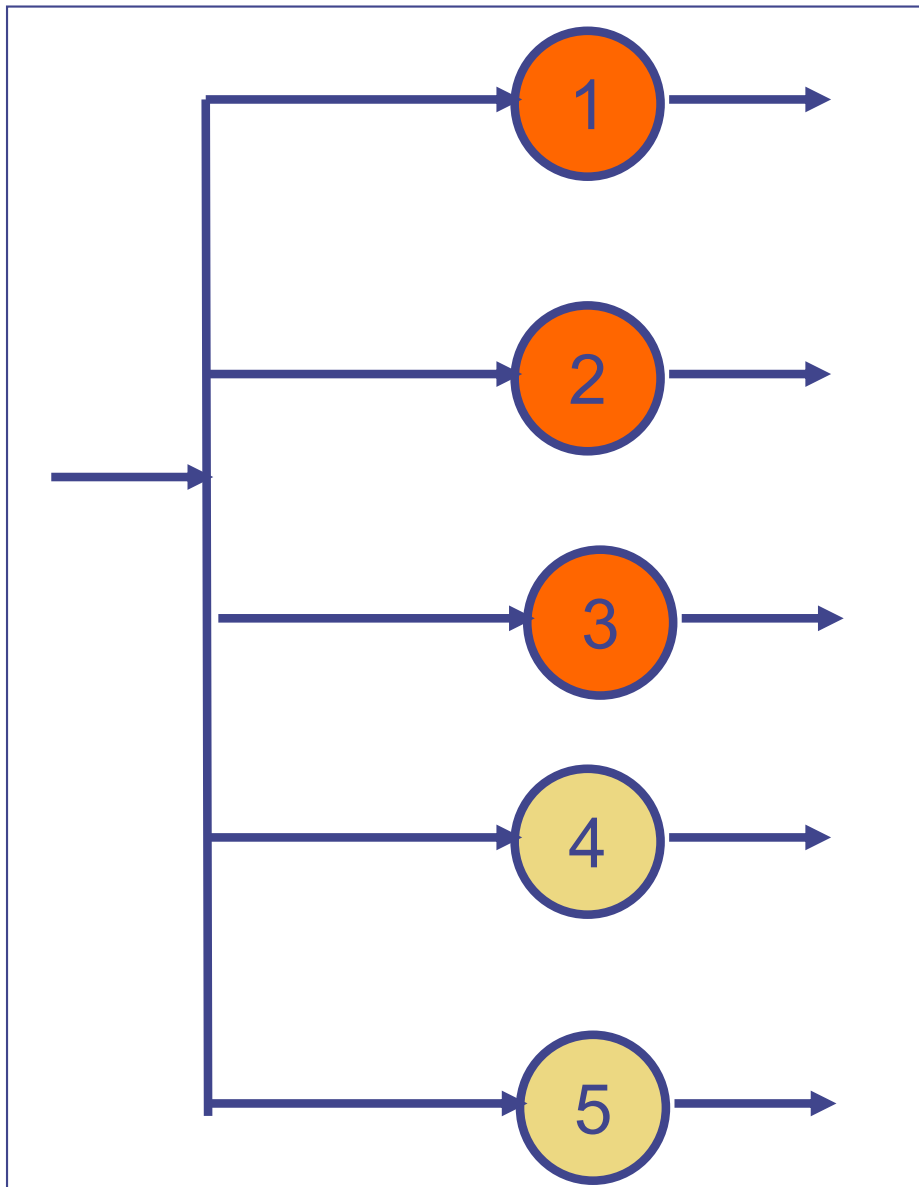
^a Computer Science Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA

^b Department of Mathematics and Computer Science, Eindhoven University of Technology, 5600 MB Eindhoven, The Netherlands

Server farm with setup cost: The problem

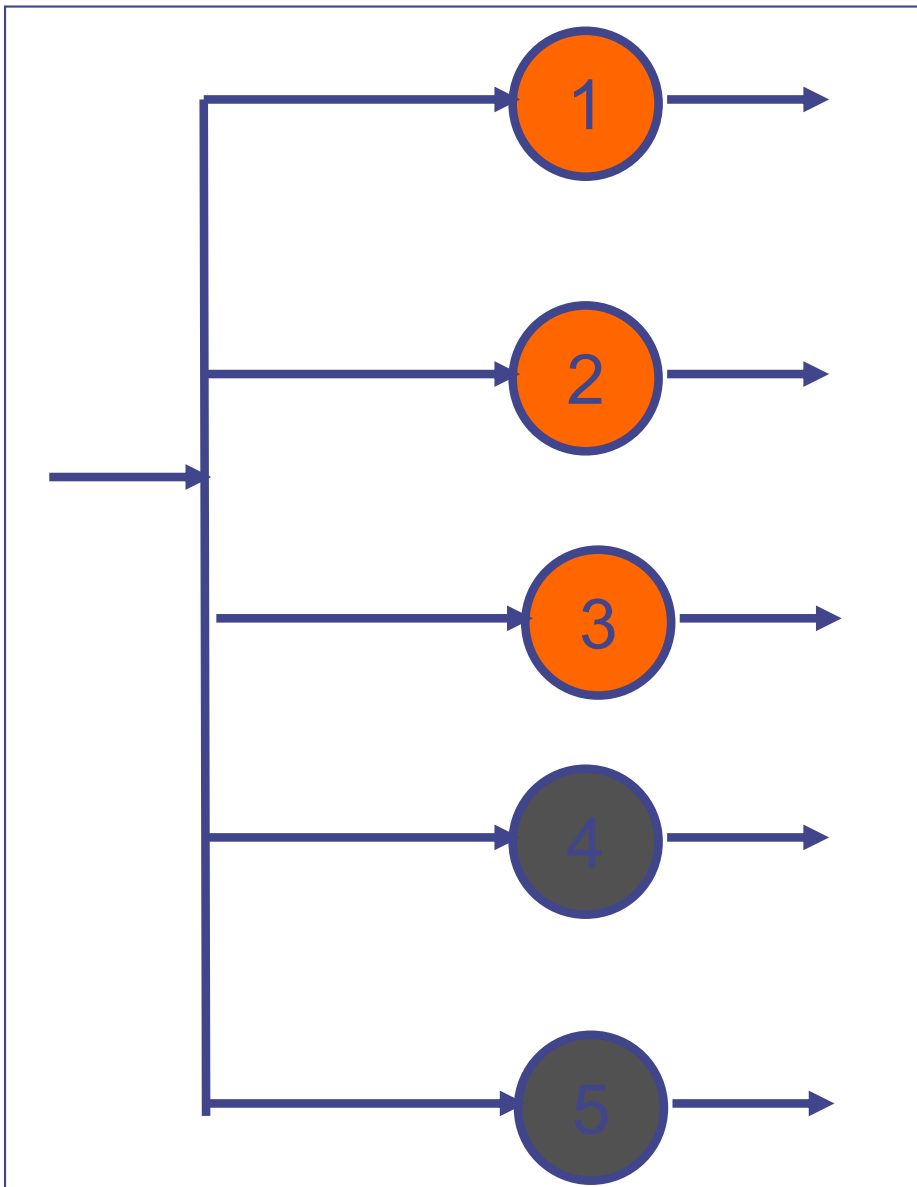
- A data centre consists of many servers
- An ON server consumes peak power
- An idling server consumes 60% of peak power
- Not all servers are needed all the time

Alternative 1: Idling servers remain ON all the time



- Servers 1-3 are ON
- Servers 4-5 are idling
- When a new job arrives, it is allocated to an idling server immediately
 - Little time delay
 - Idling servers consume power

Alternative 2: Turn idling server OFF



- Servers 1-3 are ON
- Servers 4-5 are OFF
- Once a server becomes idle, turn it off immediately to save power.
- When a new job arrives, need to turn an OFF server ON
 - Longer time delay
 - Consumes less power

Trading off between power consumption and delay

- Alternative 1: Keep idling server ON
 - Short delay but high power consumption
- Alternative 2: Turn idling server OFF immediately
 - Long delay but lower power consumption
- Can you suggest some alternatives?
- Given what you have learnt so far, will you carry out a study on this problem?

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