Introductory Queueing Theory Tutorial

Mor Harchol-Balter Computer Science Dept, Carnegie Mellon Univ.



Outline

I. Basic Vocabulary

- \circ $% \lambda$ Avg arrival rate, λ
- \circ Avg service rate, μ
- $\circ \quad \text{Avg load, } \rho$
- \circ $% \left(Avg \right) = Avg throughput, X$

- Response time, T
- Little's Law
- Exponential vs. Pareto/Heavy-tailed
- Poisson Process

II. Single-server queues

- M/G/1 response time
- Inspection Paradox
- Effect of job size variability
- Effect of load

- Scheduling: FCFS, PS, SJF, LAS, SRPT
- Scheduling: Priority Classes
- Scheduling: SOAP Framework (New)

III. Multi-server queues

- Single shared queue, M/G/k
- Load balancing across queues
- Cycle stealing

- Replication of jobs (New)
- Multi-task jobs and fork-join (New)
- Networks of queues

Vocabulary





S = job size = service requirement

$$E[S] = \frac{1}{\mu} \sec$$

Example:

- On average, job needs 3x10⁶ cycles
- Server executes 9x10⁶ cycles/sec

- Avg service rate $\mu = 3 \frac{\text{jobs}}{\text{sec}}$
- Avg size of job on this server: $E[S] = \frac{1}{3}$ sec.

Vocabulary: Load



 ρ = Load (utilization) = Frac. time server busy = $\lambda E[S] = \frac{\lambda}{\mu}$

Example:

•
$$\lambda = 2 \frac{\text{jobs}}{\text{sec}}$$
 arrive

• Each job requires $E[S] = \frac{1}{3}$ sec on avg

$$\rho = \frac{2}{3}$$

Vocabulary: Throughput

<u>Defn</u>: Throughput X is the average rate at which jobs complete (jobs/sec)



Vocabulary: Throughput



$$X=\lambda$$
 (assuming no jobs dropped)

Vocabulary: Response Time



T = response time

- T_Q = queueing time (waiting time)
- N = Number jobs in system

Little's Law:
$$E[T] = \frac{E[N]}{\lambda}$$

Vocabulary: Response Time



T = response time

$$T_Q =$$
 queueing time (waiting time)

- Q: Given that $\lambda < \mu$, what causes wait?
- A: Variability in the arrival process & service requirements

Vocabulary: Response Time







11

- "Memoryless"
- Lower variability
- Light-tail: top 1% of jobs comprise 5% load.



Decreasing hazard rate Infinite variance Heavy-tail: top 1% of jobs comprise 50% load. $S \sim \text{Pareto}(\alpha = 1)$ $\Pr\{S > x\} = \frac{1}{x}$ <u>1</u> 2 14 456 23 89 7

- "Memoryless"
- Lower variability
- Light-tail: top 1% of jobs comprise 5% load.





Variability



Vocabulary: Poisson Process with rate λ



(Poisson process comes up when aggregating many users)

Outline

I. Basic Vocabulary

- \circ $% \lambda$ Avg arrival rate, λ
- \circ Avg service rate, μ
- Avg load, ρ
- \circ Avg throughput, X

- Response time, T
- Little's Law
- Exponential vs. Pareto/Heavy-tailed
- De Poisson Process

II. Single-server queues

- M/G/1 response time
- Inspection Paradox
- Effect of job size variability
- Effect of load

- Scheduling: FCFS, PS, SJF, LAS, SRPT
- Scheduling: Priority Classes
- Scheduling: SOAP Framework (New)

III. Multi-server queues

- Single shared queue, M/G/k
- Load balancing across queues
- Cycle stealing

- Replication of jobs (New)
- Multi-task jobs and fork-join (New)
- Networks of queues

Single-Server Queue



Q: Does low $\rho \rightarrow \text{low } E[T_{Q}]$?

Single-Server Queue



A: low load does NOT ensure low wait

M/G/1



A: low load does NOT ensure low wait

Waiting for the bus



Waiting for the bus

S: time between buses

 $E[S] = 10 \min$



Waiting for the bus



M/G/1



To drop load, we can increase server speed.

Q: What can we do to combat job size variability?A: Smarter scheduling!

Scheduling in M/G/1



Well-studied scheduling policies:

- FCFS (First-Come-First-Served, non-preemptive)
- **PS** (Processor-Sharing, preemptive)
- **SJF** (Shortest-Job-First, non-preemptive)
- SRPT (Shortest-Remaining-Processing-Time, preemptive)
- LAS (Least-Attained-Service First, preemptive)

Scheduling in M/G/1



FCFS (First-Come-First-Served, non-preemptive)
PS (Processor-Sharing, preemptive)
SJF (Shortest-Job-First, non-preemptive)
SRPT (Shortest-Remaining-Processing-Time, preemptive)
LAS (Least-Attained-Service First, preemptive)

Priority Classes



According to Ruth Williams (genetic networks):

- Jobs \rightarrow molecules
- Server \rightarrow enzyme
- Classes \rightarrow protein species
- Reneging → dilution
- Class 1's load and variability can really affect class 2

Big Scheduling Breakthrough

[Scully, Harchol-Balter, Scheller-Wolf SIGMETRICS 2018]

The SOAP framework:



Enables first analysis of many previously intractable policies:

- SERPT: Prioritize jobs by Expected Remaining Size
- Gittins: Prioritize jobs by their Gittins Index
- Discretized Policies: Preemptions only at specific ages
- Mixed Priority Classes: Priority classes, where each class can have its own scheduling policy.

Outline

I. Basic Vocabulary

- \circ $% \lambda$ Avg arrival rate, λ
- \circ Avg service rate, μ
- Avg load, ρ
- \circ Avg throughput, X

- Besponse time, T
- \circ Waiting time, T_O
- Exponential vs. Pareto/Heavy-tailed
- Poisson Process

II. Single-server queues

- M/G/1 response time
- Inspection Paradox
- Effect of job size variability
- Effect of load

- Scheduling: FCFS, PS, SJF, LAS, SRPT
- Scheduling: Priority Classes
- Scheduling: SOAP Framework (New)

III. Multi-server queues

- Single shared queue, M/G/k
- Load balancing across queues
- Cycle stealing

- Replication of jobs (New)
- Multi-task jobs and fork-join (New)
- Networks of queues

M/G/k Model



Q: How does M/G/k compare with M/G/1 at k-speed?

A: Both worse and better!

Load Balancing Model



Load Balancing Model



Smart Load Balancing -> Much reduced mean response time

Cycle Stealing Model (N-model)



[Gardner, Harchol-Balter, Scheller-Wolf Transactions on Networking 2017] [Gardner, Harchol-Balter, Scheller-Wolf Operations Research 2017]

[Gardner, Harchol-Balter, Scheller-Wolf Transactions on Networking 2017] [Gardner, Harchol-Balter, Scheller-Wolf Operations Research 2017]



Same job goes to multiple queues. Job is "done" as soon as first copy completes.

[Gardner, Harchol-Balter, Scheller-Wolf Transactions on Networking 2017] [Gardner, Harchol-Balter, Scheller-Wolf Operations Research 2017]



Same job goes to multiple queues. Job is "done" as soon as first copy completes.

[Gardner, Harchol-Balter, Scheller-Wolf Transactions on Networking 2017] [Gardner, Harchol-Balter, Scheller-Wolf Operations Research 2017]



Same job goes to multiple queues. Job is "done" as soon as first copy completes.

[Gardner, Harchol-Balter, Scheller-Wolf Transactions on Networking 2017] [Gardner, Harchol-Balter, Scheller-Wolf Operations Research 2017]



Replication Tradeoff:

- + Lower response time because only need first completion.
- + Higher response time due to extra load.



















"Limited Fork-Join" See [Wang, Harchol-Balter, Jiang, Scheller-Wolf, Srikant, 2018].

Networks of Queues Model



Conclusion

I. Basic Vocabulary

- \circ $% \lambda$ Avg arrival rate, λ
- \circ Avg service rate, μ
- \circ Avg load, ρ
- \circ $% \left(Avg \right) = Avg throughput, X$

- Response time, T
- Little's Law
- Exponential vs. Pareto/Heavy-tailed
- Poisson Process

II. Single-server queues

- M/G/1 response time
- Inspection Paradox
- Effect of job size variability
- Effect of load

- Scheduling: FCFS, PS, SJF, LAS, SRPT
- Scheduling: Priority Classes
- Scheduling: SOAP Framework (New)

III. Multi-server queues

- Single shared queue, M/G/k
- Load balancing across queues
- Cycle stealing

- Replication of jobs (New)
- Multi-task jobs and fork-join (New)
- Network of queues

THANK YOU!

www.cs.cmu.edu/~harchol/

