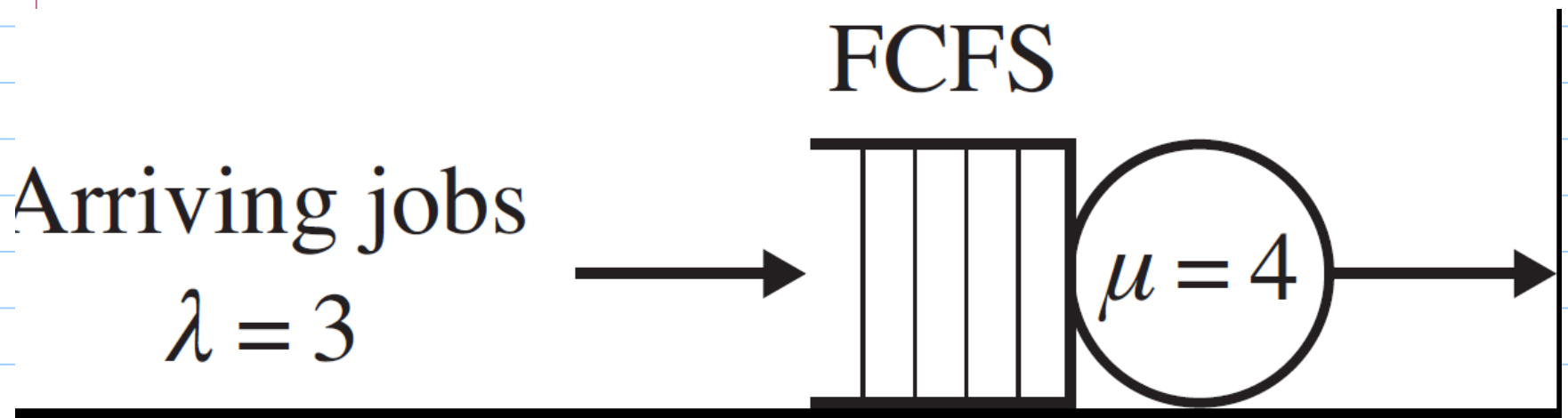


317 2015-01-15

Note Title

2015-01-15



Single-server network.

Parameters associated with the above

- Service order

By default, FCFS

(First come first Served)

- Average Arrival Rate, λ , in jobs/sec

- Mean Interarrival Time, $\frac{1}{\lambda}$, avg time between successive job arrivals

- Service Requirement Size, S , in seconds.

The time needed to process a job in the absence of queuing

- Mean Service Time, the expected value of S , $E[S]$

- Average Service Rate, $\mu = \frac{1}{E[S]}$, in jobs/sec.

Some performance metrics

T , Response Time or Turnaround Time or Time in System or Sojourn Time

$$T = t_{\text{depart}} - T_{\text{arrive}}$$

$$E[T] = E[T_Q] + E[S]$$

time in queue or delay
waiting time

$\text{Var}[T]$ (variance of Response Time)

$P\{T > t\}$ tail behavior of T

Number of Jobs in the System, N
(includes the jobs in queue and the
job being served, if any)

Number of Jobs in Queue, N_Q

We assume that $\lambda < \mu$
(average arrival rate is less than the average service
rate)

Suppose that the arrival and service distributions
are deterministic (both constant).

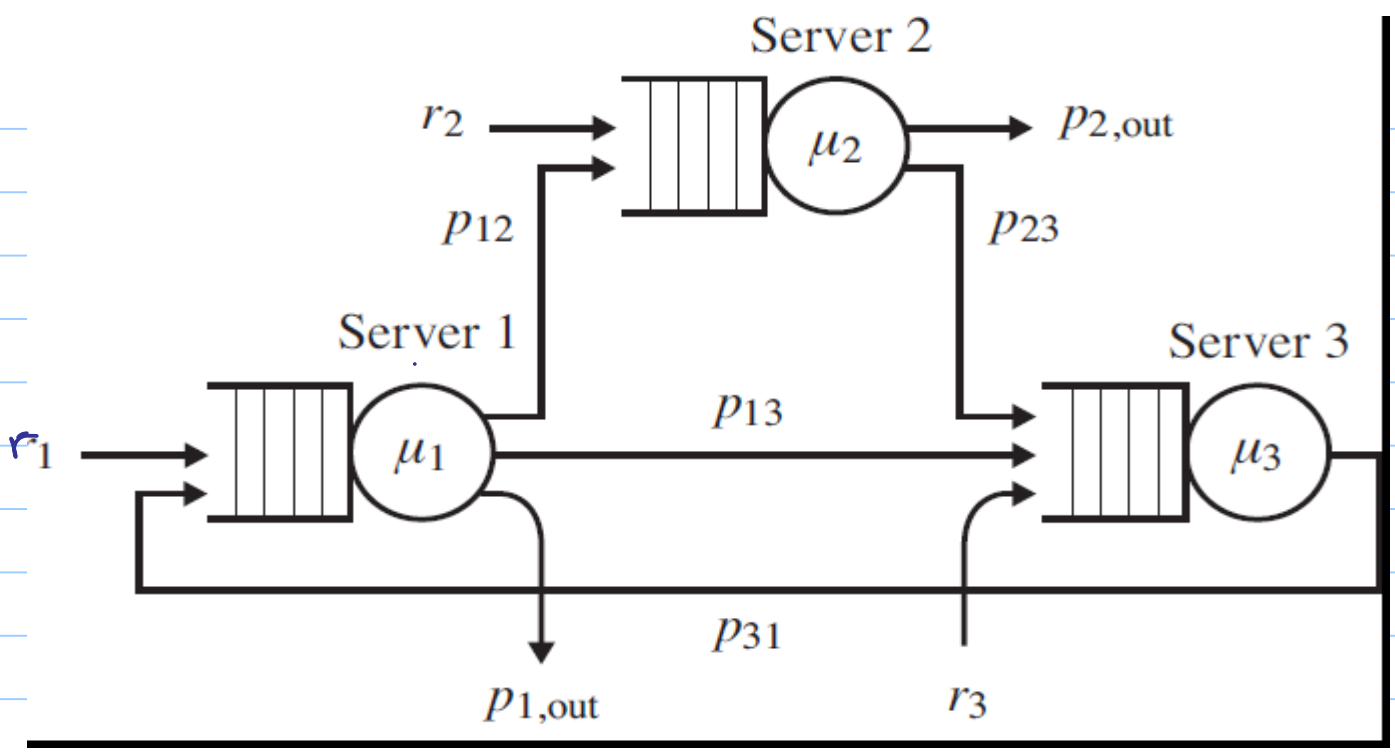
What is T_a ? What is T ?

Answer. , $T_a = 0$, $T = S$

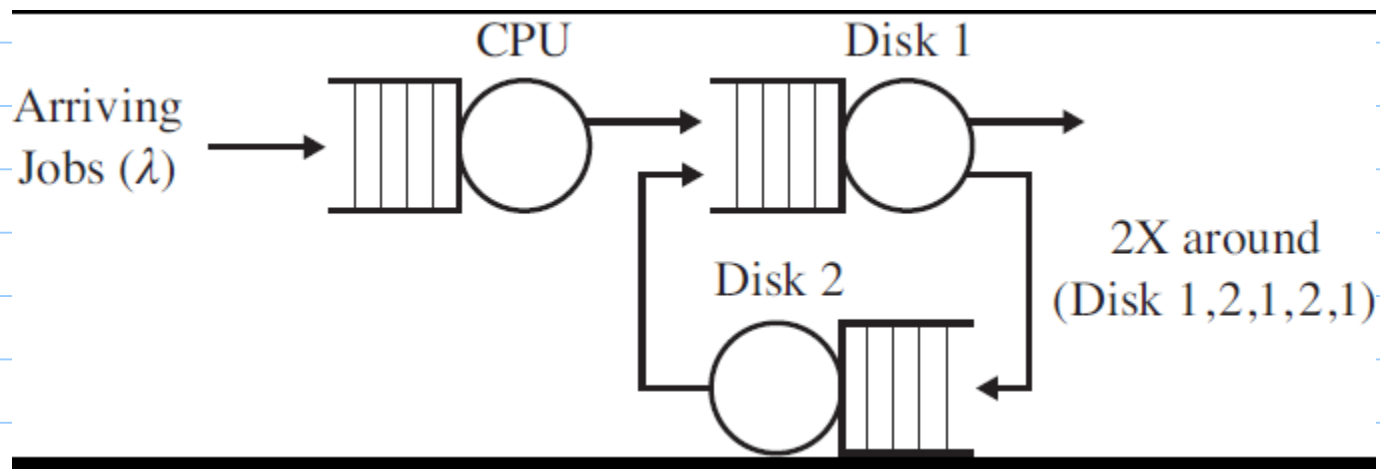
Open networks

1. Single server (see above)

2. Network of Queues w/ Probabilistic Routing



3. Network of Queues with non-Probabilistic Routing



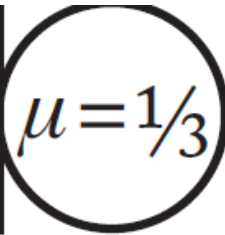
4. Finite Buffer



Space for 9 jobs
plus 1 in service

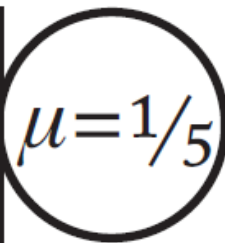
(2.5) More Metrics: Throughput & Utilization

$$\lambda = 1/6$$



versus

$$\lambda = 1/6$$



Which of these systems has the higher throughput?

Device

• Utilization ρ_i is the fraction of time device i is busy.

$$\rho_i = \frac{B}{\tau}, \text{ where } \tau \text{ is a long time interval and } B \text{ is the time within } \tau \text{ that device } i \text{ is busy.}$$

Device Throughput X_i is the rate of completions at device i (jobs/sec)

$$X_i = \frac{C}{\tau}, \text{ where } \tau \text{ is a long time interval and } C \text{ is the \# jobs completed during } \tau.$$

$$X_i = \frac{C}{r} = \frac{C}{B} \cdot \frac{B}{\alpha} = M_i \cdot \rho_i = \frac{1}{E[S_i]} \cdot \rho_i$$

\parallel \parallel
 avg ρ_i
 service
 rate
 M_i

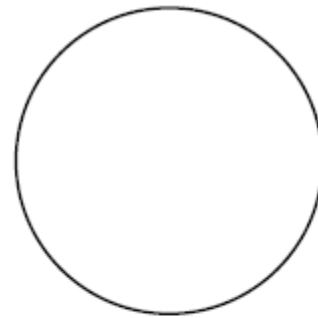
$E[S_i]$
 mean
 service time

$$\rho_i = X_i \cdot E[S_i] \quad (\text{Utilization Law})$$

$$\frac{\text{sec}}{\text{sec}} = \frac{\text{jobs}}{\text{sec}} \cdot \frac{\text{sec}}{\text{jobs}}$$

What is the throughput here \downarrow ?

$$\lambda = 1/6$$



$$\mu = 1/3$$

$X = e \cdot \mu$. We have μ . What is e ?

e = fraction of time server is busy :

$$\begin{aligned} &= \frac{\text{average service time required by a job}}{\text{average time between arrivals}} = \\ &= \frac{1/\mu}{1/\lambda} = \frac{\lambda}{\mu} = \frac{\text{average arrival rate}}{\text{average service rate}} \end{aligned}$$

[Warning: not a formal proof.]

Proof is given on p. 100 - Corollary 6.5]

Substitute $e = \frac{\lambda}{\mu}$ in $X = e \cdot \mu$, and obtain

$$X = \frac{\lambda}{\mu} \cdot \mu = \lambda \quad (! !)$$

The throughput does not depend on the service rate!