



Lecture 10

Prof. Krishna R. Pattipati

**Dept. of Electrical and Computer Engineering
University of Connecticut**

Contact: krishna@engr.uconn.edu (860) 486-2890

ECE 336



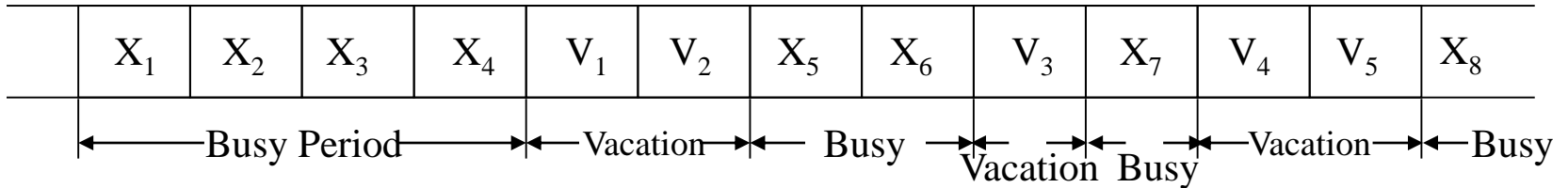
Outline of Lecture 10

- M|G|1 Queues with vacations
- Application of M|G|1 Results to Reservations & Polling
- Application to Token Ring Networks
- Extension to Non-product form Networks with M|G|1 Nodes
- M|G|1 Queues with Priorities
- Extensions to Queuing Networks with Priorities
- G|G|1 Queues



Queues with Vacations

- Suppose that at the end of a busy period, the server goes on “vacation” for some random interval of time. Thus, a new arrival to an idle system, rather than going into service immediately, waits for the end of a vacation period.



Usual $M/G/1$: Alternate **busy-idle-busy** cycles
 $M/G/1$ with vacations: busy - vacation - - - busy
can be multiple vacation cycles

- Variations:**
 - The server may continue taking vacations until, on return from a vacation, it finds at least one customer ... multiple vacations model
 - The server takes exactly one vacation. Single vacation model
 - Busy-vacation-idle-busy-vacation-busy- ...cycles



Application of Queues with Vacations -1

Applications:

1. Machine breakdown:
 - Serving Customers (Primary customers)
 - Maintaining the system when machine fails (secondary customers)

Can also be viewed as a priority model:

- Two priorities
- Breakdowns have preemptive priority over the primary customers.

2. Maintenance in production systems:

- During idle periods, we do preventive maintenance on the system. The system is assumed to never breakdown during production (i.e., busy period)
- This is a single vacation model.



Application of Queues with Vacations -2

3. Maintenance in computer and communication systems:

- Processor in computer and communication systems perform considerable testing and maintenance to improve reliability. There exist several ways in which the maintenance functions are scheduled in these situations:
- Maintenance is divided into short segments. Whenever the primary jobs are absent, the processor does a segment of the maintenance work. If upon completion of the segment, there are no primary jobs, the system continues with the next segment of maintenance. This corresponds to “multiple vacation models” case
- There exist variations on this model. Some examples include:
 - Always make sure that a certain amount of time δ is spent on maintenance
 - For every m primary jobs, do one maintenance job. This is a **limited service vacation model**, in which the server takes vacation on becoming idle or after serving m consecutive primary jobs.
 - For every T seconds on primary jobs, spend on one segment on the maintenance job.
 - Periodic maintenance \Rightarrow secondary jobs with preemptive or non preemptive priority over primary jobs

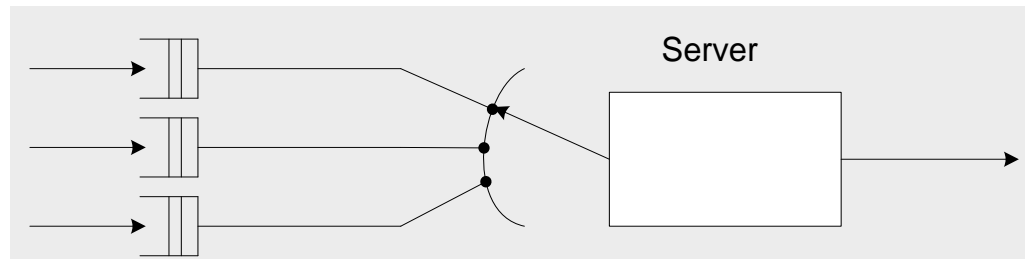


Application of Queues with Vacations - 3

4. Cyclic server queues:

- These arise naturally as models of schedules in computer systems and communication networks, e.g., task processing in computer systems, scheduling virtual circuits or ports in a communication system.
- The basic model here has m classes of customers, each with its own queue
- These m queues are served by a single server cyclically.

Question: When does the server move from one queue to the next?



- Exhaustive Service: The server leaves a queue when it is empty
⇒ Multiple vacation model
- Gated Service: The server upon arrival to a queue, closes a gate behind the waiting customers in that queue, and leaves that queue when the customers present before the gate is closed are served.



Application of Queues with Vacations - 4

- **Limited Service:** There is a limit R_i placed on the number of customers served on each visit to queue i . The server leaves queue i when that queue is empty or when R_i customers have been served during the current visit.

■ We will consider $M/G/1$ queue with vacations and applications of cyclic server queues to communication networks.

- Poisson arrival process
- V_1, V_2, V_3, \dots are i.i.d random variables.
- Service times are i.i.d random variables.

■ $M/G/1$ multiple vacation case:

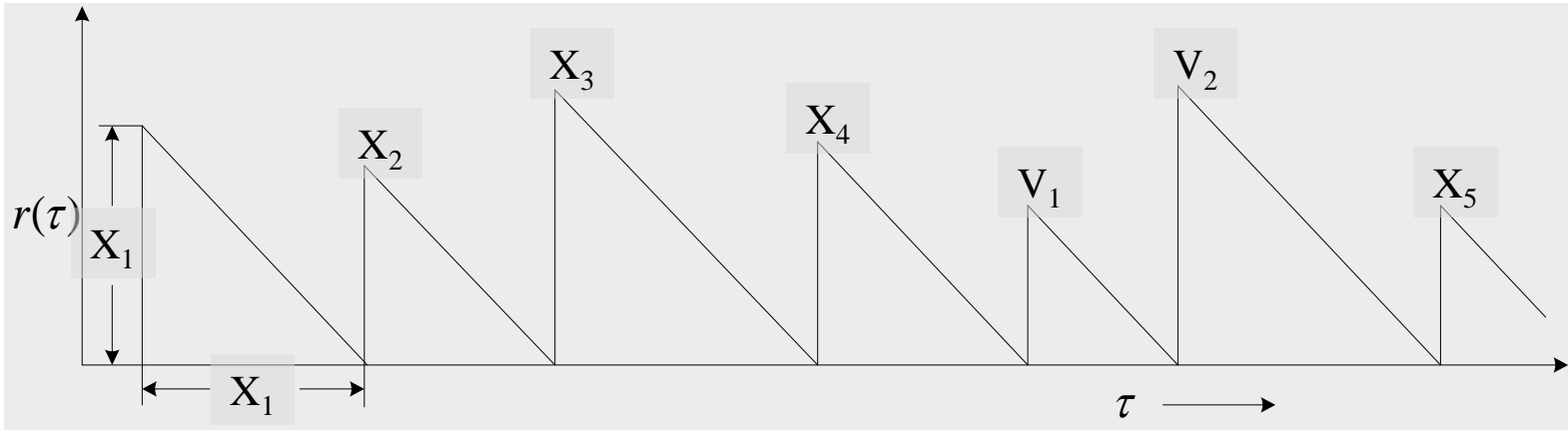
- **What does a new arrival do?**
 - Wait in the queue for the completion of current service and then the service of *all* customers waiting before it.
 - Wait for vacation

$$\Rightarrow W = \frac{X_R}{1 - \rho}$$

X_R = mean residual time of completion of service or vacation in process when the i^{th} customer arrives.



M|G|1 Queue with Multiple Vacations - 1



Time averages = Ensemble Averages

$$\frac{1}{t} \int_0^t r(\tau) d\tau = \frac{1}{t} \left[\frac{1}{2} \sum_{i=1}^{M(t)} X_i^2 + \frac{1}{2} \sum_{i=1}^{L(t)} V_i^2 \right]$$

$M(t)$ = number of service completions in $(0, t)$

$L(t)$ = number of vacations in $(0, t)$

As $t \rightarrow \infty$,

$$\begin{aligned} \overline{X_R} &= \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t r(\tau) d\tau \\ &= \frac{1}{2} \left[\lim_{t \rightarrow \infty} \frac{M(t)}{t} \cdot \frac{1}{M(t)} \sum_{i=1}^{M(t)} X_i^2 + \frac{1}{2} \lim_{t \rightarrow \infty} \frac{L(t)}{t} \cdot \frac{1}{L(t)} \sum_{i=1}^{L(t)} V_i^2 \right] \end{aligned}$$



M|G|1 Queue with Multiple Vacations - 2

As $t \rightarrow \infty$, fraction of time occupied with vacations is $(1 - \rho)$

$$\text{Total vacation time} = (1 - \rho)t$$

$$\text{Average vacation time } \bar{V} = \frac{(1 - \rho)t}{L(t)}$$

or

$$\boxed{\frac{L(t)}{t} = \frac{1 - \rho}{\bar{V}}}$$

Also, $\lambda = \frac{M(t)}{t}$

So, $\bar{X}_R = \frac{1}{2} \lambda \bar{X}^2 + \frac{1}{2} (1 - \rho) \frac{\bar{V}^2}{\bar{V}}$

$$\begin{aligned} W_{M/G/1/V_M} &= \frac{1}{2} \frac{\lambda \bar{X}^2}{(1 - \rho)} + \frac{1}{2} \frac{\bar{V}^2}{\bar{V}} \\ &= W_{M/G/1} + \text{Residual Vacation Time} \end{aligned}$$

Indeed, This decomposition is valid in a wider generality.

See B.T. Doshi, Journal of Applied Probability, Vol. 22, pp.419-428, 1985



M|G|1 Queue with a Single Vacation - 1

- M/G/1 queue with a single vacation: HW problem. See Doshi (1985) and Fuhrmann, Operations Research, 1984, pp.1368-1373

Result :

$$W = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \frac{\lambda \bar{V}}{B_v(\lambda) + \lambda \bar{V}} \cdot \frac{\bar{V}^2}{2\bar{V}}$$

where $B_v(\lambda) = \int_0^\infty e^{-\lambda V} f_v(V) dV$

Hint:

$$(1-\rho)t = L(t)(\bar{V} + \bar{I})$$

$$I = \begin{cases} 0 & \text{if } \tau_a < V \\ \tau_a - V & \text{if } \tau_a > V \end{cases}$$

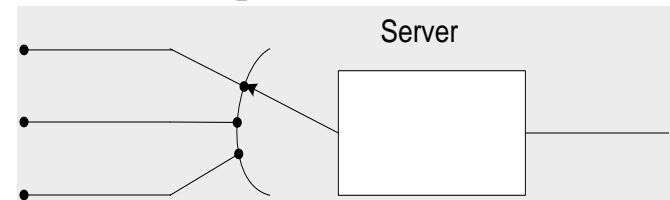
$\tau_a = \text{inter-arrival time}$

- Application to Communication Networks: Cyclic queues

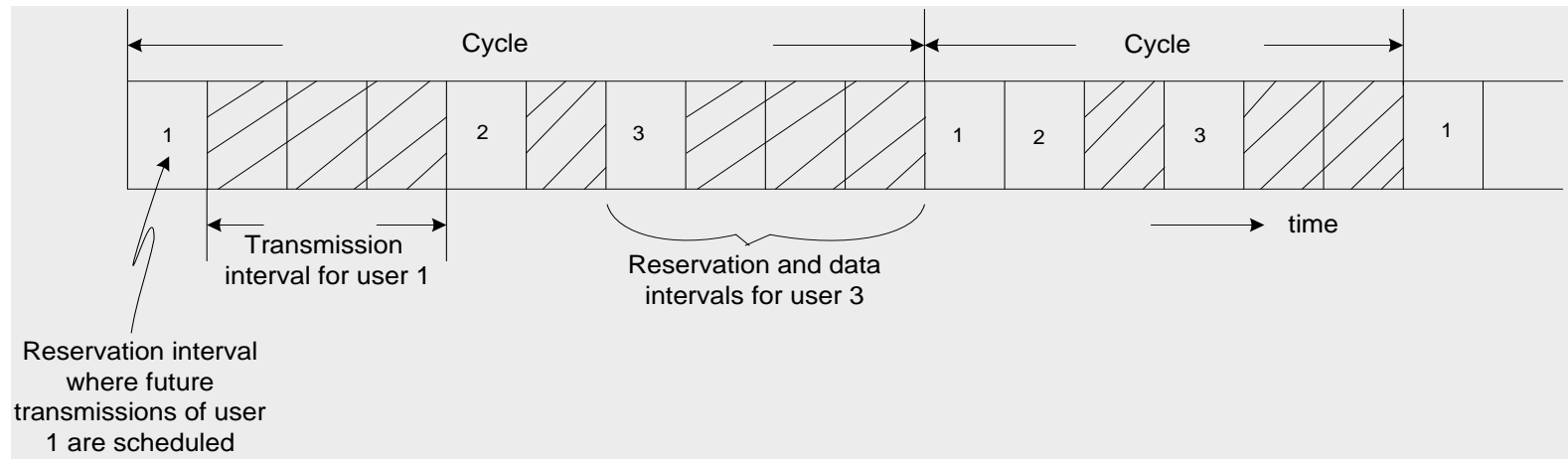
- A communication channel is accessed by several spatially separated users.
- Only one user can transmit successfully on the channel at one time
=> a multi-access channel
- Communication resource of the channel can be divided into two portions:

→ Packet transmissions ... data intervals

→ Reservation (or Polling) messages that schedule future packet Transmissions ... reservation intervals



Cyclic Queues - 1



- m users
- Assume that each data interval contains packets of a single user
- Reservations for these packets are made in the immediately preceding reservation interval
- All users are taken up in cyclic order $(1, 2, 3, \dots, m, 1, 2, 3, \dots)$
- Three versions depending on how packets are transmitted during the data interval of each user

Exhaustive system: A packet of a user that occurs during the user's reservation or data interval is transmitted in the same data interval \Rightarrow channel goes to the next user only after completing the transmission of all the packets of the current user ... Token ring

Cyclic Queues - 2

Partially gated system: Only packets that arrived until the end of the reservation interval are transmitted during the current data interval.

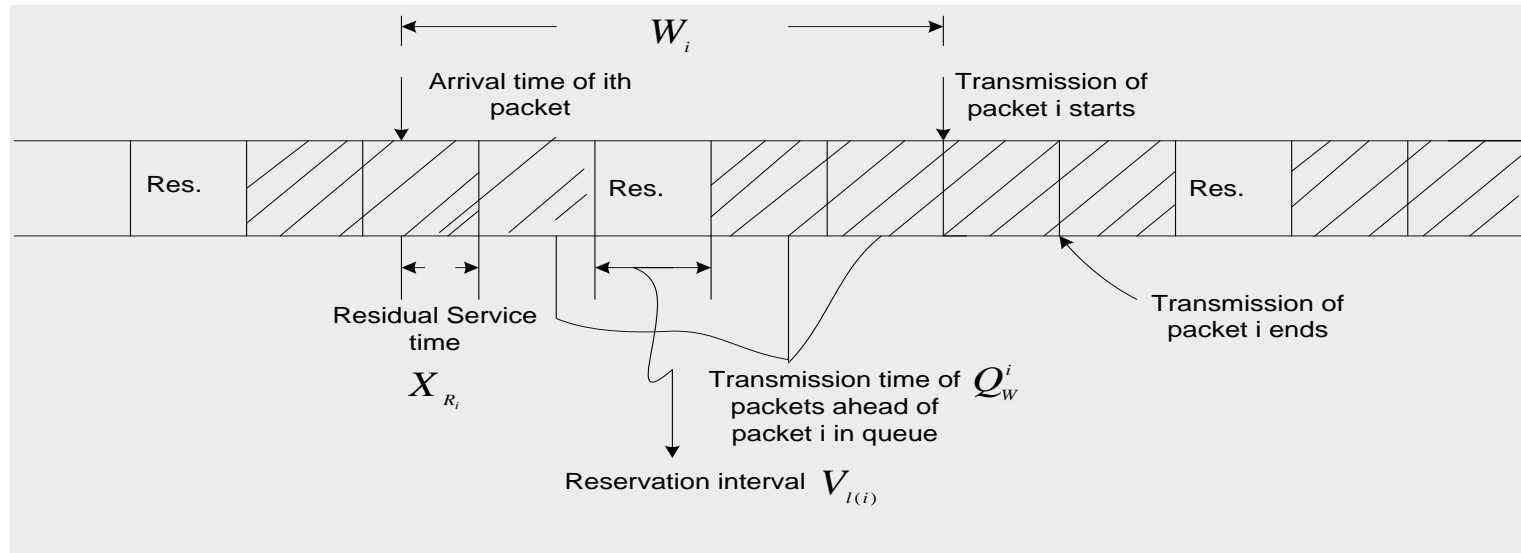
Gated system: Only packets that arrived prior to the reservation interval can be transmitted.

■ Analysis:

- Arrival processes of all users are Poisson with rate λ/m
- 1st and 2nd moments of packet transmission times $\overline{X} = \frac{1}{\mu}$ and $\overline{X^2}$ (i.i.d. random variables)
- Inter-arrival times and packet transmission times are independent
- Reservation intervals of different users can have different distributions, but we assume it to be the same for simplicity.

Cyclic Queues - 3

Single-user system: $m=1$ (Gated)



$V_l = l^{\text{th}}$ reservation interval. Successive reservation intervals are i.i.d with 1^{st} and 2^{nd} moments \bar{V} and $\overline{V^2}$

$\{V_l, X_i, \tau_a\}$ are independent. $\tau_a =$ inter-arrival time

$$E[W_i] = E[X_{R_i}] + E[Q_{\omega_i}^i] \bar{X} + E[V_{l(i)}]$$

Cyclic Queues - 4

Consider a gated system

$$X_R = \frac{\lambda \overline{X^2}}{2} + \frac{(1-\rho)\overline{V^2}}{2\overline{V}}$$

$$W = \frac{\lambda \overline{X^2}}{2} + \frac{(1-\rho)\overline{V^2}}{2\overline{V}} + \rho W + \overline{V}$$

$$\text{So, } W = \frac{\lambda \overline{X^2}}{2(1-\rho)} + \frac{\overline{V^2}}{2\overline{V}} + \frac{\overline{V}}{1-\rho} \quad (\text{single user, gated})$$

Cyclelength

Suppose $\overline{V} = A$ (deterministic), then

$$\begin{aligned} W &= \frac{\lambda \overline{X^2}}{2(1-\rho)} + \frac{A}{2} + \frac{A}{1-\rho} \\ &= W_{M/G/1} + \frac{A}{2} \left[\frac{3-\rho}{1-\rho} \right] \end{aligned}$$

Exhaustive system: $\Rightarrow M/G/1$ with vacations

$$W = \frac{\lambda \overline{X^2}}{2(1-\rho)} + \frac{\overline{V^2}}{2\overline{V}} \quad \overline{V} = A \Rightarrow W = W_{M/G/1} + \frac{A}{2}$$

Similar to $M/G/1$ with vacations
 \Rightarrow a vacation starts when all previous
arrivals are served.

Cyclic Queues - 5

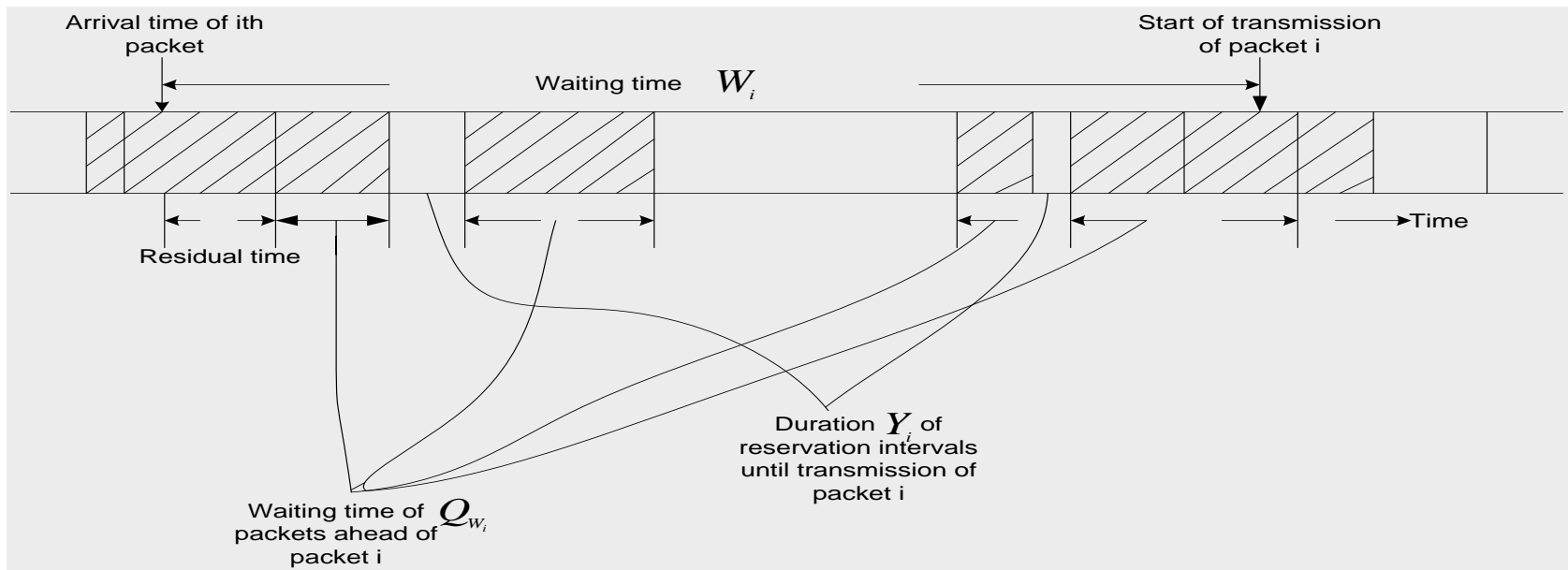
Partially gated:

$$W = X_R + \rho W + \rho \bar{V}$$

$$(or) W = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \frac{\bar{V}^2}{2\bar{V}} + \frac{\rho \bar{V}}{1-\rho}$$

$$\bar{V} = A \text{ deterministic} \Rightarrow W = W_{M/G/1} + \frac{A}{2} \left(\frac{1+\rho}{1-\rho} \right)$$

Multi-user system:



Cyclic Queues - 6

- user data rates λ/m for users $0,1,2,\dots,m-1$
- l^{th} reservation interval is used to make reservations for user $l \bmod(m) = l - \left\lfloor \frac{l}{m} \right\rfloor \cdot m$
and the subsequent l^{th} data interval is used to send packets corresponding to those reservations
- Consider packet i

$$E\{W_i\} = E\{X_{Ri}\} + E\{Q_{oi}\} \bar{X} + E\{Y_i\}$$

as $i \rightarrow \infty$

$$W = \bar{X}_R + \rho W + Y$$

$$\Rightarrow W = \frac{\bar{X}_R + Y}{1 - \rho}$$

$$\text{Know } \bar{X}_R = \frac{\lambda}{2} \bar{X}^2 + \frac{(1 - \rho) \sum_{k=0}^{m-1} \bar{V}_k^2}{2 \sum_{k=0}^{m-1} \bar{V}_k} = \frac{\lambda}{2} \bar{X}^2 + \frac{(1 - \rho) \bar{V}^2}{2 \bar{V}}$$

$$\text{Need to compute } Y : m = 1 \Rightarrow Y = \begin{cases} 0 \text{ exhaustive} \\ \rho \bar{V} \text{ partially gated} \\ \bar{V} \text{ gated} \end{cases}$$

Cyclic Queues - 7

- What happens when $m > 1$? Consider Exhaustive case

Let

$$\alpha_{ij} = E \left\{ \begin{array}{l} Y_i \mid \text{packet } i \text{ arrives in user } l\text{'s reservation or data interval} \\ \text{and belongs to user } (l+j) \bmod m \end{array} \right\}$$

$$\Rightarrow \alpha_{ij} = \begin{cases} 0; & j=0 \\ \bar{V}_{(l+1) \bmod m} + \dots + \bar{V}_{(l+j) \bmod m} & ; j=1, 2, \dots, m-1 \end{cases}$$

since packet i belongs to any user with probability $\frac{1}{m}$, we have

$$E \{ Y_i \mid \text{packet } i \text{ arrives in user } l\text{'s reservation or data interval} \}$$

$$= \sum_{j=1}^{m-1} \frac{m-j}{m} \bar{V}_{(l+j) \bmod m}$$

Finally, a packet will arrive during l 's data interval with probability $\frac{\rho}{m}$

a packet will arrive during l 's reservation interval with probability $\frac{(1-\rho)\bar{V}_l}{\sum_{k=0}^{m-1} \bar{V}_k}$

Cyclic Queues - 8

Let $i \rightarrow \infty$

$$\Rightarrow Y = \sum_{l=0}^{m-1} \left(\frac{\rho}{m} + \frac{(1-\rho)\bar{V}_l}{\sum_{k=0}^{m-1} \bar{V}_k} \right) \sum_{j=1}^{m-1} \frac{m-j}{m} \bar{V}_{(l+j) \bmod m} = \frac{\rho(m-1)\bar{V}}{2} + \frac{(1-\rho)m\bar{V}}{2} - \frac{(1-\rho)\sum_{k=0}^{m-1} \bar{V}_k^2}{2m\bar{V}}$$

$$\Rightarrow W_{exh} = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \frac{\sigma_v^2}{2\bar{V}} + \frac{m-\rho}{2} \frac{\bar{V}}{(1-\rho)}; \sigma_v^2 = \frac{\sum_{k=0}^{m-1} (\bar{V}_k^2 - \bar{V}_k^2)}{m}; \bar{V} = \frac{\sum_{k=0}^{m-1} \bar{V}_k}{m}$$

See Bertsekas & Gallager, pp. 200 for details

Partially gated system: Same as exhaustive, except that if a packet arrives during user's own data interval, it is delayed by an additional $m\bar{V}$.

This occurs with probability $\frac{\rho}{m}$

$$\Rightarrow Y_{\rho G} = Y_{exh} + \rho\bar{V}$$

$$W_{\rho G} = W_{exh} + \frac{\rho\bar{V}}{1-\rho}$$

Cyclic Queues - 9

Gated System : If a packet arrives during a user's own reservation or data interval, it is delayed by an additional $m\bar{V}$ time units. This occurs with probability $\frac{1}{m}$

$$Y_G = Y_{exh} + \bar{V} \Rightarrow W_G = W_{exh} + \frac{\bar{V}}{1-\rho}$$

suppose $\bar{V} = \frac{A}{m}$, then

$$W_{exh} = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \frac{A(1-\rho/m)}{2(1-\rho)} = W_{M/G/1} + \frac{A(1-\rho/m)}{2(1-\rho)} = W_{M/G/1/V_m} + \frac{A}{2} \cdot \frac{m-1}{m} \cdot \frac{\rho}{(1-\rho)}$$

$$W_{\rho G} = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \frac{A(1+\rho/m)}{2(1-\rho)}$$

$$W_G = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \frac{A(1+(2-\rho)/m)}{2(1-\rho)}$$

As $m \rightarrow \infty$

$$W_{exh} = W_{\rho G} = W_G = \frac{\lambda \bar{X}^2}{2(1-\rho)} + \underbrace{\frac{A}{2(1-\rho)}}_{1/2 \text{ cycle length}}$$

Cyclic Queues - 10

Limited service Systems: $k_i = 1$ case

- In each user's data interval, only the first packet of the user waiting in queue (if any) is transmitted rather than all waiting packets.
- Consider only gated and partially gated systems (exhaustive case doesn't make sense here)

As before

$$W = \bar{X}_R + \rho W + Y_L$$

What is Y_L ?

Consider partially gated system

A packet arriving during user l 's data or reservation interval will belong to any

one of the users with probability $\frac{1}{m}$. In steady state, the average number of packets

waiting in the individual queue of the user that owns the arriving packet = $\frac{\lambda W}{m}$

Cyclic Queues - 11

Each of these packets cause an extra cycle of reservations of length $m\bar{V}$. So,

$$Y_{L,\rho G} = Y_{\rho G} + \frac{\lambda W}{m} \cdot m\bar{V}$$
$$\Rightarrow W_{L,\rho G} = \frac{X_R + Y_{\rho G}}{(1 - \rho - \lambda\bar{V})}$$
$$= W_{\rho G} \cdot \frac{(1 - \rho)}{(1 - \rho - \lambda\bar{V})}$$

Gated System :

$$Y_{L,G_i} = Y_{L,\rho G_i} + m\bar{V} \text{ Prob} \left\{ \begin{array}{l} \text{packet } i \text{ arrives during the reservation interval of its owner} \\ \text{and the subsequent data interval is empty} \end{array} \right\}$$

$$\text{Prob}\{\text{packet } i \text{ arrives during the reservation interval}\} = 1 - \rho$$

$$\text{Let Prob}\{\text{reservation interval followed by an empty data interval}\} \stackrel{\Delta}{=} p$$

$$\text{Prob}\{\text{reservation interval followed by a nonempty data interval}\} = 1 - p$$

Cyclic Queues - 12

$$\therefore (1-p) \frac{\bar{X}}{\bar{V}} = \frac{\rho}{1-\rho} \Rightarrow \rho \bar{V} = (1-p) \bar{X} - (1-p) \rho \bar{X}$$

$$\Rightarrow 1-p = \frac{\lambda \bar{V}}{1-\rho} \Rightarrow p = \frac{(1-\rho - \lambda \bar{V})}{(1-\rho)} = 1 - \frac{\lambda \bar{V}}{1-\rho}$$

So,

$$Y_{L,G_i} = Y_{L,\rho G_i} + \frac{m \bar{V} (1-\rho - \lambda \bar{V}) (1-\rho)}{(1-\rho) m}$$

$$W_{L,G} = \left[\bar{X}_R + Y_{L,\rho G} + \bar{V} (1-\rho - \lambda \bar{V}) \right] / (1-\rho - \lambda \bar{V})$$
$$= W_{L,\rho G} \left(\frac{1-\rho}{1-\rho - \lambda \bar{V}} \right) + \bar{V}$$

- Note that we need $\lambda(\bar{X} + \bar{V}) < 1$ for stability



Application to Token Ring Networks

Application to Token ring networks:

- m poisson streams with rate $\frac{\lambda}{m}$
- \bar{V} propagation delay + relaying delay per step $\Rightarrow \bar{V} = \frac{A}{m}$
- $\bar{X} = 1 \Rightarrow \rho = \lambda$

Exhaustive System:

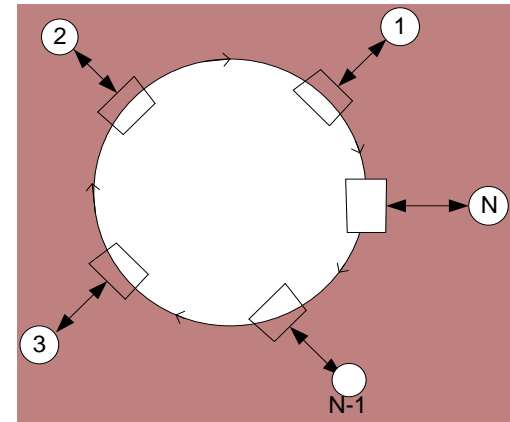
$$W = \frac{\lambda \bar{X}^2}{2(1-\lambda)} + \frac{A}{2} \frac{m-\rho}{m(1-\rho)} = W_{M/G/1} + \frac{\bar{V}(m-\lambda)}{2(1-\lambda)}$$

$$\bar{V} = \frac{A}{m}$$

Partially gated, limited service system:

$$W = \frac{\lambda \bar{X}^2 + (m + \lambda) \bar{V}}{2(1 - \lambda - \lambda \bar{V})}$$

$$\Rightarrow \text{stable if } \lambda < \frac{1}{1 + \bar{V}}$$



TDM vs. FDM vs. SFDM

- Slotted time-division and frequency-division multiplexing:

- FDM:
- m poisson streams with rate $\frac{\lambda}{m}$
 - transmission time of each packet = m time units.

$$\text{Each channel is an M/D/1 queue} \Rightarrow W_{FDM} = \frac{\lambda m}{2(1-\lambda)}$$

Slotted FDM: Packet transmissions can start only at times $0, m, 2m, \dots$ M/D/1 queue with vacations

where $\bar{V} = m, \bar{V}^2 = m^2$

$$W_{SFDM} = W_{FDM} + \frac{m}{2} = \frac{m}{2(1-\lambda)}$$

$$\text{TDM: } W_{TDM} = W_{SFDM} = W_{FDM} + \frac{m}{2} = \frac{m}{2(1-\lambda)}$$

Response Times :

$$R_{FDM} = m + \frac{\lambda m}{2(1-\lambda)}$$

$$R_{SFDM} = R_{FDM} + \frac{m}{2}$$

$$R_{TDM} = 1 + \frac{m}{2(1-\lambda)} = R_{FDM} - \left(\frac{m}{2} - 1 \right)$$

$$\underline{R_{TDM} < R_{FDM} < R_{SFDM} \text{ for } m > 1}$$



Networks with General Service Times

No product form \Rightarrow Approximation

MVA is ideally suited for these approximations

For product-form networks, have

$$R_{ij}(\underline{n}) = \frac{S_{ij}}{\mu_i} \left[1 + Q_i(\underline{n} - \underline{e}_j) \right]$$

At FCFS nodes, need $S_{ij} = S_i \quad \forall j$

1) Suppose S_{ij} is different for different classes j & exponential

$$\Rightarrow R_{ij}(\underline{n}) \cong \frac{S_{ij}}{\mu_i} + \left[\sum_{k=1}^J Q_{ik}(\underline{n} - \underline{e}_j) S_{ik} \right] \cdot \frac{1}{\mu_i}$$

2) General service demand requirement:

$$\Rightarrow R_{ij}(\underline{n}) \cong \frac{S_{ij}}{\mu_i} + \sum_{k=1}^J \left[Q_{ik}(\underline{n} - \underline{e}_j) - u_{ik}(\underline{n} - \underline{e}_j) \right] \frac{S_{ik}}{\mu_i} + \sum_{k=1}^J u_{ik}(\underline{n} - \underline{e}_j) \frac{\tilde{S}_{ik}}{\mu_i}$$

$$\tilde{S}_{ik} = \text{residual service demand} = \frac{S_{ik}^2 + \sigma_{ik}^2}{2S_{ik}}$$



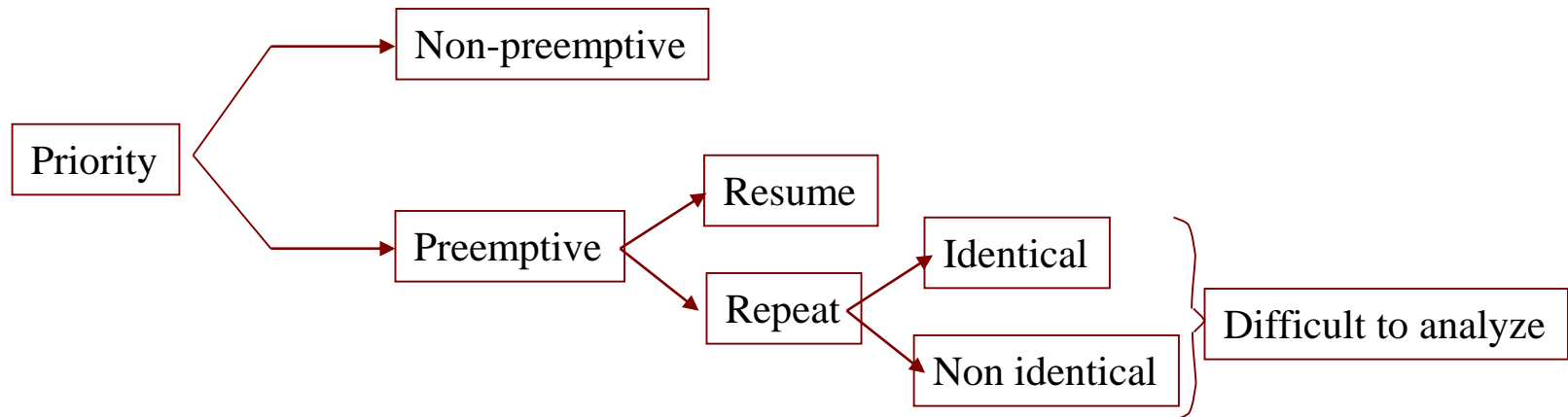
M|G|1 Queues with Priorities -1

References:

1. Kleinrock, Vol.II, Chapter 3
2. N.K.Jaiswal, "Priority Queues", Academic Press, 1968

- Customers are divided into J different priority classes
Class 1: higher priority Class J : least priority
Static priorities predetermined (not dependent on waiting time, # in the system, etc.)
- Arrival processes are independent, Poisson & independent of the service times

$$\lambda_k, \quad \bar{X}_k = \frac{1}{\mu_k}, \quad \bar{X}_k^2$$





M|G|1 Queues with Priorities - 2

■ **Nonpreemptive priority:** a customer is allowed to complete service without interruption even if a customer of higher priority arrives in the mean time

- A separate queue for each priority class
- When the server becomes free, the first customer of the highest nonempty priority queue enters service => Head-Of- the-Line (HOL) priority.
- Need to compute waiting time of each priority class. We appeal to conceptual reasoning rather than analytic derivation

Q_{Wk} : Average number waiting in queue k

W_k : Average waiting time for priority k customers

$$\rho_k = \frac{\lambda_k}{\mu_k} \quad \text{System Utilization for priority } k \text{ customers}$$

\overline{X}_R : Mean residual time

Assume

$$\rho_1 + \rho_2 + \dots + \rho_J < 1$$

If not $\exists a k^* \ni W_{k^*+i} = \infty \quad i = 0, 1, 2, \dots, J - k^*$

$$W_1 = \overline{X}_R + Q_{W1} \overline{X}_1 \Rightarrow W_1 = \frac{\overline{X}_R}{1 - \rho_1}$$



M|G|1 Queues with Priorities - 3

$$W_2 = \bar{X}_R + Q_{W_1} \bar{X}_1 + Q_{W_2} \bar{X}_2 + \overbrace{\lambda_1 W_2 \bar{X}_1}^{\text{Future arrivals}}$$

$$= \bar{X}_R + \rho_1 W_1 + \rho_2 W_2 + \rho_1 W_2 \Rightarrow W_2 = \frac{\bar{X}_R + \rho_1 W_1}{1 - \rho_2 - \rho_1} = \frac{\bar{X}_R}{(1 - \rho_1)(1 - \rho_1 - \rho_2)}$$

In general,

$$W_k = \bar{X}_R + \rho_1 W_1 + \rho_2 W_2 + \dots + \rho_k W_k + \left(\sum_{i=1}^{k-1} \rho_i \right) W_k$$

$$\Rightarrow W_k = \frac{\bar{X}_R}{\left[1 - \sum_{i=1}^{k-1} \rho_i \right] \left[1 - \sum_{i=1}^k \rho_i \right]}$$

$$R_k = W_k + \bar{X}_k$$

Need \bar{X}_R :

$$\bar{X}_R = \frac{1}{2} \left(\sum_{i=1}^J \lambda_i \right) \bar{X}^2$$

$$\bar{X}^2 = \frac{\lambda_1}{\sum_{i=1}^J \lambda_i} \bar{X}_1^2 + \frac{\lambda_2}{\sum_{i=1}^J \lambda_i} \bar{X}_2^2 + \dots + \frac{\lambda_J}{\sum_{i=1}^J \lambda_i} \bar{X}_J^2$$

$$\Rightarrow \bar{X}_R = \frac{1}{2} \sum_{i=1}^J \lambda_i \bar{X}_i^2$$

Waiting time of a high priority class (e.g., W_1) depends on the arrival rates of lower priority classes



M|G|1 Queues with Priorities - 4

Note:

$$1) W_k = \bar{X}_R + \sum_{l=1}^k \rho_l W_l + \left[\sum_{l=1}^{k-1} \rho_l \right] W_k$$

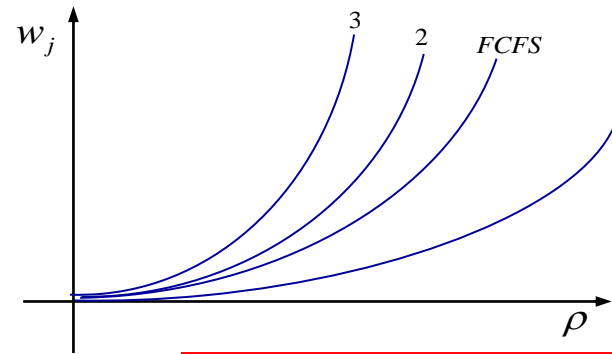
$$\sum_{k=1}^J \rho_k W_k = \rho \bar{X}_R + \sum_{k=1}^J \sum_{l=1}^k \rho_l \rho_k W_l + \sum_{k=1}^J \sum_{l=1}^{k-1} \rho_l \rho_k W_k$$

$$\sum_{k=1}^J \rho_k W_k = \rho \bar{X}_R + \sum_{l=1}^J \left(\sum_{k=l}^J \rho_k \right) \rho_l W_l + \sum_{k=1}^J \sum_{l=1}^{k-1} \rho_l \rho_k W_k$$

$$= \rho \bar{X}_R + \sum_{k=1}^J \left(\sum_{l=1}^J \rho_l \right) \rho_k W_k$$

$$\Rightarrow \sum_{k=1}^J \rho_k W_k = \frac{\rho \bar{X}_R}{1 - \rho} \quad \text{independent of priorities}$$

M/G/1 Conservation law



Some customer classes do better than FCFS, Others do worse.

For multi-server queues, see Buzen, Operations Research, 1983

2) Suppose it costs $\$C_k$ per unit of wait in the queue by a class k customer

Want to minimize

$$\sum_{k=1}^J C_k Q_{Wk} = \sum_{k=1}^J \frac{C_k}{X_k} (\rho_k W_k)$$

$$\Rightarrow \min \sum_{k=1}^J \left(\frac{C_k}{X_k} \right) \rho_k W_k$$

The μC Rule

Suppose classes are ordered according to [1] [2] ... [k][k+1] [k+2]...[J]... Optimal \Rightarrow Cost = $\sum_{k=1}^J \frac{C_{[k]}}{X_{[k]}} \rho_{[k]} W_{[k]}$

Suppose we interchange [k] and [k+1]

Does not affect waiting times of [1] [2] ... [k-1] & [k+2] [k+3]...[J]

$$\Rightarrow \frac{C_{[k]}}{X_{[k]}} \rho_{[k]} W_{[k]}^* + \frac{C_{[k+1]}}{X_{[k+1]}} W_{[k+1]}^* \rho_{[k+1]} \leq \frac{C_{[k+1]}}{X_{[k+1]}} W_{[k+1]} \rho_{[k+1]} + \frac{C_{[k]}}{X_{[k]}} W_{[k]} \rho_{[k]}$$

$$\frac{C_{[k+1]}}{X_{[k+1]}} \rho_{[k+1]} \Delta W_{[k+1]} \leq -\frac{C_{[k]}}{X_{[k]}} \rho_{[k]} \Delta W_{[k]}; \Delta W_{[k]} = W_{[k]}^* - W_{[k]} < 0; \Delta W_{[k+1]} = W_{[k+1]}^* - W_{[k+1]} > 0$$

$$\Rightarrow \left[\frac{C_{[k+1]}}{X_{[k+1]}} - \frac{C_{[k]}}{X_{[k]}} \right] \rho_{[k+1]} \Delta W_{[k+1]} \leq 0 \quad \text{since} \quad \sum_{k=1}^J \rho_{[k]} \Delta W_{[k]} = 0$$

$$\Rightarrow \frac{C_{[k+1]}}{X_{[k+1]}} \leq \frac{C_{[k]}}{X_{[k]}} \quad \text{since} \quad \Delta W_{[k+1]} > 0 \Rightarrow C_{[k+1]} \mu_{[k+1]} \leq C_{[k]} \mu_{[k]}$$

$$\text{or} \quad C_{[k+1]} X_{[k]} \leq C_{[k]} X_{[k+1]}$$

$$\text{or} \quad \frac{X_{[k]}}{C_{[k]}} \leq \frac{X_{[k+1]}}{C_{[k+1]}}$$

arrange priorities according to μC rule or weighted shortest-processing time rule (WSPT)

$$\mu_1 C_1 \geq \mu_2 C_2 \geq \dots \geq \mu_J C_J$$

- μC rule minimizes expected waiting cost



M/G/1 with Preempt-resume Priority - 1

KEY: The waiting time of a high priority customer class is *independent of* the arrival rates of lower priority classes (unlike non-preemptive priority)

- * service of a lower priority customer is interrupted when a high priority customer arrives, and is resumed from the point of interruption once all customers of higher priority have been served

Here, we find it convenient to calculate the response time rather than the waiting time. Consider class j . R_j consists of :

$$R_j = Term_a + Term_b + Term_c$$

Term a: Average service time $\bar{x}_j = 1/\mu_j$ (since preempt-resume)

Term b: Average time required, upon arrival of a priority j customer, to service customers of priorities 1 to j , that are already in the system \Rightarrow average unfinished work corresponding to priorities 1 through j .



M/G/1 with Preempt-resume Priority - 2

What is Term b?

The average waiting time of an $M/G/1$ queue with arrivals due to classes $1, 2, \dots, j$ (priorities $j+1, j+2, \dots, J$ are neglected)

$$Term_b = \frac{\sum_{i=1}^j \lambda_i \bar{x}_i^2}{2(1 - \sum_{i=1}^j \rho_i)} = \frac{\bar{x}_{Rj}}{(1 - \sum_{i=1}^j \rho_i)}$$

Term c: Average waiting time for customers of priorities 1 through $(j-1)$ who arrive while the customer of class j is in the system

What is Term c?

$$Term_c = \sum_{i=1}^{j-1} \frac{1}{\mu_i} \lambda_i R_j = \left(\sum_{i=1}^{j-1} \rho_i \right) R_j$$

$$\therefore R_j = \frac{1}{\mu_j} + \frac{\bar{x}_{Rj}}{(1 - \sum_{i=1}^j \rho_i)} + \left(\sum_{i=1}^{j-1} \rho_i \right) R_j \Rightarrow R_j = \frac{\frac{1}{\mu_j} (1 - \sum_{i=1}^j \rho_i) + \bar{x}_{Rj}}{(1 - \sum_{i=1}^j \rho_i)(1 - \sum_{i=1}^{j-1} \rho_i)}$$



M/G/1 with Preempt-resume Priority - 3

- Can recursively evaluate R_j :

```
RHOSUM=0
 $\bar{x}_R = 0$ 
Do  $j = 1, \dots, J$ 
  TEMP=RHOSUM
  RHOSUM = RHOSUM +  $\rho_j$ 
   $\bar{x}_R = \bar{x}_R + \frac{\lambda_j x_j^2}{2}$ 
   $R_j = \frac{\frac{1}{\mu_j} (1 - RHOSUM) + \bar{x}_R}{(1 - RHOSUM)(1 - TEMP)}$ 
End do
```

- Extension to multiple servers: Buzen, *Operations Research*, 1983
Agrawal *Metamodeling*, MIT Press 1985



Extension to Multi-class Queuing Networks - 1

- Two approximations:
 - Shadow approximation, due to *Sevcik*, valid for only preempt-resume discipline
 - *Bryant, Lakshmi, Chandy* and *Krzenski* approximation (also termed MVA approximation) ***The Best***

Assume only single server nodes. Infinite server nodes are easy; multi-server & state-dependent nodes are research issues.

- M nodes $\{1, 2, \dots, M\}$
- J Classes $\{1, 2, \dots, J\}$, class 1 has highest priority, ..., class J has lowest
- visits, v_{ij} ;
- Mean service time per visit : $s_{ij} = 1/\mu_{ij}$
- R_{ij} = Response time over all visits;
- Q_{ij} = Queue length at node i for class j ;
- $X_{\bar{j}}$ = Throughput of class j customers



Extension to Multi-class Queuing Networks - 2

■ Shadow approximation:

K. Sevcik, "Priority scheduling disciplines in QN models of computing systems", Proc. IFIP congress, North Holland, 1977, pp. 565-574

■ Preempt resume service discipline. Assume a single PR center

■ Key idea:

1. Replace each priority center by J shadow centers, where J is the number of priority classes
2. *Each shadow service center is visited by one class only*
3. Service time per visit of class j customers at the shadow service center is

$$s_{sj} = \frac{s_{pj}}{1 - \sum_{k=1}^{j-1} \rho_{pk}}; \quad \rho_{pk} = x_k s_{pk}$$

Solve $(M+J-1)$ node product-form network



Extension to Multi-class Queuing Networks - 3

- Algorithm:

Initialize ρ_{pk}

While ρ_{pk} not converged Do

$$s_{sj} = \frac{S_{pj}}{1 - \sum_{k=1}^{j-1} \rho_{pk}}$$

Solve $(M+J-1)$ product form network $\Rightarrow x_j$

Compute ρ_{pk}

End

- Can easily extend to multiple *preempt-resume* service centers



Extension to Multi-class Queuing Networks - 4

Algorithm:

Initialize ρ_{pk} at all $p \in P_R$

While ρ_{pk} , $p \in P_R$ not converged Do

$$S_{sj} = \frac{S_{pj}}{1 - \sum_{k=1}^{j-1} \rho_{pk}} ;$$

Solve ($|P_R| J + M - |P_R|$) node
product form network

Compute ρ_{pk}

End

Errors can be as high as 40%!!



Extension to Multi-class Queuing Networks - 5

- MVA approximation:

Bryant, M.S.Lakshmi, K.M.Chandy and A.E.Krzenski “MVA Priority Approximation”, ACM Trans. on Comp. Systems, Feb. 1983

- Recall product-form MVA equations

Repeat $\forall \underline{n} \in 0 \leq \underline{n} \leq N$

$$R_{ij}(\underline{n}) = v_{ij} s_{ij} [1 + Q_i(\underline{n} - \underline{e}_j)]$$

$$x_j(\underline{n}) = n_j / \sum_{i=1}^M R_{ij}(\underline{n})$$

$$Q_{ij}(\underline{n}) = x_j(\underline{n}) R_{ij}(\underline{n})$$

End Loop

- Restrictions:

$$QD \sim PS \text{ or } LCFS PR$$

$$QD \sim FCFS \Rightarrow s_{ij} = s_i \text{ independent of } j$$

How do we extend these results to queuing networks with priority nodes?



Extension to Multi-class Queuing Networks - 6

- Suppose we have an isolated open node with arrival rates $\lambda_1, \lambda_2, \dots, \lambda_j$ with visits $v_j = 1$. Then,
 - **Preempt-resume**

$$R_j = s_j + \sum_{k=1}^j (Q_k - \rho_k) s_k + \sum_{k=1}^{j-1} R_j \lambda_k s_k + \bar{s}_{Rj};$$

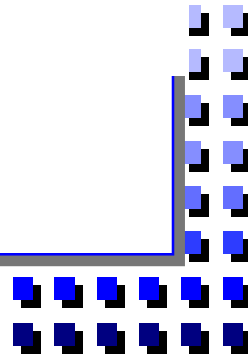
$$\bar{s}_{Rj} = \frac{\sum_{k=1}^j \lambda_k \bar{s}_k^2}{2} = \sum_{k=1}^j \frac{\rho_k s_k [1 + C_{s_k}^2]}{2}; \quad C_{s_k} = \frac{\sigma_{s_k}}{s_k}$$

$$\Rightarrow R_j = \frac{\left[s_j + \sum_{k=1}^j \left\{ (Q_k - \rho_k) s_k + \frac{\lambda_k \bar{s}_k^2}{2} \right\} \right]}{1 - \sum_{k=1}^{j-1} \rho_k} \quad *$$

$$= \frac{\left[s_j + \sum_{k=1}^j \left\{ Q_k s_k + \frac{\rho_k s_k (C_{s_k}^2 - 1)}{2} \right\} \right]}{1 - \sum_{k=1}^{j-1} \rho_k}$$

For exponential case:

$$R_j = \frac{\left[s_j + \sum_{k=1}^j Q_k s_k \right]}{1 - \sum_{k=1}^{j-1} \rho_k}$$



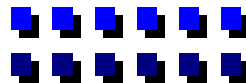


Extension to Multi-class Queuing Networks - 7

– Non-preemptive:

$$\begin{aligned}
 W_j &= \underbrace{\sum_{k=1}^j (Q_k - \rho_k) s_k}_{\text{waiting time due to customers ahead of our tagged customer}} + \underbrace{\overline{s_R}}_{\text{remaining service time}} + \underbrace{\sum_{k=1}^{j-1} W_j \lambda_k s_k}_{\text{waiting time due to arrivals after our tagged customer came in}} \\
 &= \sum_{k=1}^j (Q_k - \rho_k) s_k + \sum_{k=1}^J \frac{\rho_k s_k [1 + C_{sk}^2]}{2} + W_j \sum_{k=1}^{j-1} \rho_k \\
 \Rightarrow R_j &= W_j + s_j = s_j + \frac{\sum_{k=1}^j (Q_k - \rho_k) s_k + \sum_{k=1}^J \frac{\rho_k s_k [1 + C_{sk}^2]}{2}}{1 - \sum_{k=1}^{j-1} \rho_k} \quad ** \\
 \text{For Exponential case: } R_j &= s_j + \frac{\sum_{k=1}^j Q_k s_k + \sum_{k=j+1}^J \rho_k s_k}{1 - \sum_{k=1}^{j-1} \rho_k}
 \end{aligned}$$

Equation (*) and (**) form the basis of MVA equations for priority networks





Extension to Multi-class Queuing Networks - 8

- Consider preempt-resume priority case first
 - Assumption 1: Poisson arrival s at the nodes ----- not true in networks

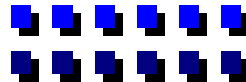
$$R_{ij}(\underline{n}) = \frac{v_{ij} [s_{ij} + \sum_{k=1}^j Q_{ik}^{(j)}(\underline{N}) s_{ik}] + \sum_{k=1}^j \frac{\rho_{ik}^{(j)} s_{ik} [c_{sk}^2 - 1]}{2}}{[1 - \sum_{k=1}^{j-1} \rho_{ik}^{(j)}(\underline{N})]}$$

- $Q_{ik}^{(j)}(\underline{N})$ = average number of class k customers at node i as seen by an arrival of a customer of class j ... $Q_{ik}(\underline{N} - \underline{e}_j)$ for product-form networks
- Assumption 2: we assume that the arrival theorem is valid

$$\Rightarrow Q_{ik}^{(j)} = Q_{ik}(\underline{N} - \underline{e}_j)$$

$$R_{ij}(\underline{n}) = \frac{v_{ij} [s_{ij} + \sum_{k=1}^j Q_{ik}(\underline{N} - \underline{e}_j) s_{ik}] + \sum_{k=1}^j \frac{\rho_{ik}^{(j)} s_{ik} [c_{sk}^2 - 1]}{2}}{[1 - \sum_{k=1}^{j-1} \rho_{ik}^{(j)}(\underline{N})]}$$

- Need means of computing ρ_{ik}





Extension to Multi-class Queuing Networks - 9

– Assumption 3: Know

$$\rho_{ik}^{(j)}(\underline{N}) = x_k^{(j)}(\underline{N})s_{ik}(\underline{N})v_{ik}$$

$$\rho_{ik}^{(j)}(\underline{N}) = \rho_{ik}(\underline{N} - \underline{e}_j) \quad \text{Not Good !!}$$

$$\rho_{ik}^{(j)}(\underline{N}) = \rho_{ik}(\underline{N}) \quad \text{Not good when utilization of server} > 0.7$$

(Bryant–Krzenski approximation)

$$\rho_{ik}^{(j)}(\underline{N}) = \rho_{ik}(\underline{N} - Q_{ik}\underline{e}_k) \quad \text{Best approximation !}$$

(Chandy-Lakshmi approximation)

When there are Q_{ik} customers at node i , the arrival rate of class k at node i is determined by the $(N_k - Q_{ik})$ customers in the network

$$\Rightarrow \rho_{ik}(\underline{N} - Q_{ik}\underline{e}_k) = x_{ik}(\underline{N} - Q_{ik}\underline{e}_k)s_{ik}v_{ik}$$

– Errors generally less than 10%. Extension to non-preemptive is easy (**)

- Open problems:
 - Method validated for exponential service times. General service times open.
 - Extension to multi-server & state-dependent server modes
 - B-S and C-N approximations for priority MVA

M/G/1 Busy periods

$$P_0 = \text{Prob idle} = \lim_{n \rightarrow \infty} \frac{I_1 + I_2 + \cdots + I_n}{(I_1 + I_2 + \cdots + I_n) + (B_1 + B_2 + \cdots + B_n)} = \frac{E(I)}{E(I) + E(B)}$$

- For M/G/1

$$E(I) = \frac{1}{\lambda}$$

Also, know $1 - P_0 = \lambda E(x) = \lambda \bar{x} = \rho$

$$\Rightarrow 1 - \lambda \bar{x} = \frac{1/\lambda}{1/\lambda + \bar{B}} \quad \Rightarrow \quad \bar{B} = \frac{\bar{x}}{1 - \lambda \bar{x}}$$

Average # of customs served per busy period = $\frac{1}{1 - \rho}$

M/G/1 with Batch Arrivals

Batch arrivals:

$$\alpha_j = \text{Prob}\{\text{batch size} = j\}$$

$$\text{Expected Batch Size} = E(N) = \sum_{j=0}^{\infty} j\alpha_j$$

$$W = \bar{X}_R + \lambda E(N) \bar{X} W \Rightarrow W = \frac{\bar{X}_R}{1 - \lambda E(N) \bar{X}}$$

$$\begin{aligned} \bar{X}_R = \text{remaining service time of customer in service} \\ + \text{waiting time due to those in batch} \end{aligned} = \frac{\lambda E(X^2) E(N)}{2} + E(W_B)$$

$$E(W_B) = \sum_j E(W_B | \text{batch size} = j) \text{Prob}\{\text{batch size} = j\}$$

$$\text{Prob}\{\text{batch size} = j\} = \frac{j\alpha_j}{\sum_j j\alpha_j} = \frac{j\alpha_j}{E(N)}$$

$$E(W_B | \text{batch size} = j) = \sum_{i=1}^j (i-1) \bar{X} \frac{1}{j} = \frac{j-1}{2} \bar{X} \Rightarrow E(W_B) = \sum_j (j-1) \alpha_j j \frac{\bar{X}}{2E(N)} = \frac{\bar{X}[E(N^2) - E(N)]}{2E(N)}$$

$$\text{so, } W = \frac{\left\{ \frac{\lambda E(X^2) E(N)}{2} + \frac{\bar{X}[E(N^2) - E(N)]}{2E(N)} \right\}}{1 - \lambda E(N) \bar{X}}$$



G/G/1 Waiting Time Bound - 1

- Can get only bounds on waiting time

$$W \leq \frac{\lambda(\sigma_a^2 + \sigma_x^2)}{2(1 - \rho)} \quad \text{equality as } \rho \rightarrow 1$$

σ_a^2 : variance of inter-arrival times

σ_x^2 : variance of service times

λ : arrival rate

$$\rho = \frac{\lambda}{\mu}$$

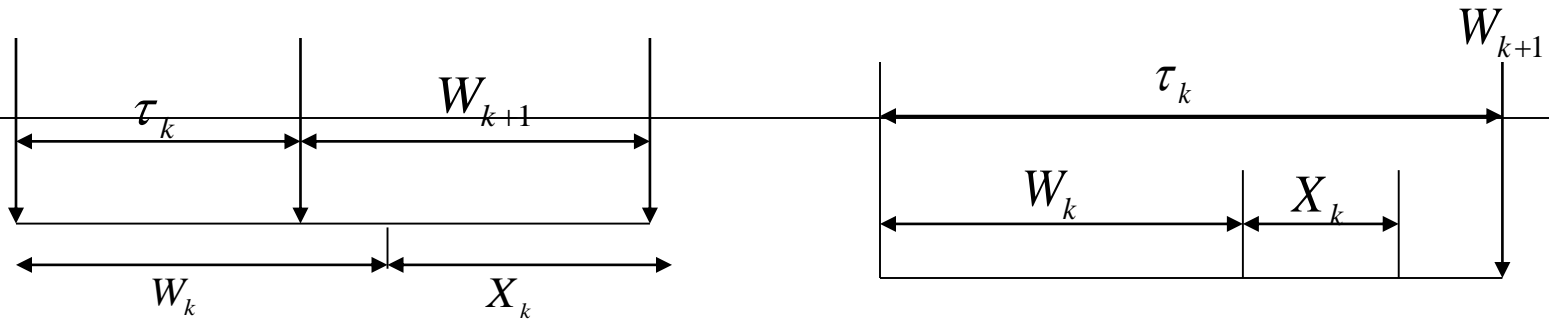
Let W_k : waiting time of k^{th} customer

X_k : service time of k^{th} customer

τ_k : inter-arrival time between k^{th} and $(k+1)^{\text{th}}$ customer



G/G/1 Waiting Time Bound - 2



$$W_{k+1} = \max(0, W_k + X_k - \tau_k)$$

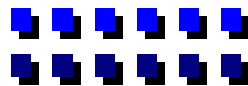
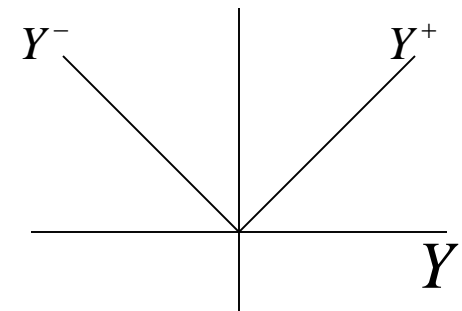
$$= \max(0, W_k + V_k)$$

■ Some identities:

$$Y^+ = \max(0, Y) \qquad Y^- = -\min(0, Y)$$

$$\Rightarrow Y = Y^+ - Y^- \quad \& \quad Y^+ Y^- = 0$$

$$\bar{Y} = \bar{Y}^+ - \bar{Y}^- \qquad \sigma_Y^2 = \sigma_{Y^+}^2 + \sigma_{Y^-}^2 + 2\bar{Y}^+ \bar{Y}^-$$





G/G/1 Waiting Time Bound - 3

$$W_{k+1} = (W_k + V_k)^+ ; V_k = X_k - \tau_k ; I_k = (W_k + V_k)^-$$

$$\begin{aligned} \sigma_{(W_k+V_k)}^2 &= \sigma_{(W_k+V_k)^+}^2 + \sigma_{(W_k+V_k)^-}^2 + 2\overline{(W_k + V_k)^+} \overline{(W_k + V_k)^-} \\ &= \sigma_{W_k}^2 + \sigma_{V_k}^2 = \sigma_{W_k}^2 + \sigma_a^2 + \sigma_X^2 \quad (\text{since } W_k \text{ and } V_k \text{ are independent}) \end{aligned}$$

$$\sigma_{W_k}^2 + \sigma_a^2 + \sigma_X^2 = \sigma_{W_{k+1}}^2 + \sigma_{I_k}^2 + 2\overline{W}_{k+1} \overline{I}_{k+1}$$

$$\text{As } k \rightarrow \infty, \overline{W}_k \rightarrow \overline{W}; \overline{I}_k \rightarrow \overline{I}; \sigma_{W_k}^2 \rightarrow \sigma_W^2$$

$$\Rightarrow \quad \overline{W} = \frac{\sigma_a^2 + \sigma_X^2}{2\overline{I}} - \frac{\sigma_I^2}{2\overline{I}}$$

$$\text{Average idle time } \overline{I} = \frac{(1-\rho)}{\lambda} \Rightarrow \quad \overline{W} \leq \frac{(\sigma_a^2 + \sigma_X^2)\lambda}{2(1-\rho)}$$

as $\rho \rightarrow 1, \sigma_I^2 \rightarrow 0$ since $I \rightarrow 0$ with probability 1

■ Special case: M/G/1

$$W = \frac{(\sigma_a^2 + \sigma_X^2 - \sigma_I^2)}{2\overline{I}} = \frac{\lambda(\sigma_X^2 + \mu^{-2})}{2(1-\rho)}$$

$$\Rightarrow \sigma_I^2 = \frac{1}{\lambda^2} - \frac{1}{\mu^2}$$

$$\text{Neglected item in the bound: } \frac{\lambda\sigma_I^2}{2(1-\rho)} = \frac{\lambda\mu(\mu^2 - \lambda^2)}{2\lambda^2\mu^2(\mu - \lambda)} = \frac{1}{2} \left(\frac{1}{\lambda} + \frac{1}{\mu} \right) < \frac{1}{\lambda} \text{ for } \rho < 1$$



Summary

- M|G|1 Queues with vacations
- Application of M|G|1 Results to Reservations & Polling
- Application to Token Ring Networks
- Extension to Non-product form Networks with M|G|1 Nodes
- M|G|1 Queues with Priorities
- Extensions to Queuing Networks with Priorities
- G|G|1 Queues