# Priority Operation

## of

## The M/G/1 Queue

1

---

*Class*



*Class 1 Lowest Priority Class*

Head of Line (HOL) Priority Operation of M/G/1 Queue

2

```
                    ┌──────────────────────┐
                    │  Priority Discipline │
                    └──────────────────────┘
                               │
              ┌────────────────┴────────────────┐
       ┌──────────────┐                  ┌────────────┐
       │ Non-Preemptive│                  │ Preemptive │
       └──────────────┘                  └────────────┘
              ▲                                  │
              ┊                         ┌────────┴────────┐
              ┊                  ┌────────────┐    ┌─────────────┐
   work conserving              │ Preemptive │    │ Preemptive  │
     disciplines                │  Resume    │    │ Non-Resume  │
              ┊                  └────────────┘    └─────────────┘
              ┊                         │
              └┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┘
```

3

---

**Non-Preemptive Priority**

•Consider an arrival of priority class $j$ when the server is serving a job of lower priority class $k$, $j>k$.

•The new arrival, in spite of being of a priority level higher than the current job in service, will not interrupt the on-going service.

•Instead, it will join the queue (FCFS) at the end of the queue of its own priority class, i.e. Class $j$, and wait for the current job to finish service.

•Normal HOL priority operation will resume once the on-going service is over

> *On-going service is not interrupted, even if there are new arrivals of higher priority*
>
> ***Work-conserving Discipline***

4

**Preemptive Resume Priority**

•Consider an arrival of priority class $j$ when the server is serving a job of lower priority class $k$, $j>k$.

•The new arrival of class $j$ will immediately preempt the lower priority job currently being served and will start its own service.

•When service to the previously preempted class $k$ job eventually resumes (possibly after service to the preempting job of class $j$ and other jobs of priority higher than $k$), *the service is resumed from the point where it was interrupted earlier.*

> *On-going service interrupted by arrival of higher priority.*
>
> *Work already done for the preempted job is remembered*
>
> ***Work-conserving Discipline***

5

---

**Preemptive Non-Resume Priority**

•Consider an arrival of priority class $j$ when the server is serving a job of lower priority class $k$, $j>k$.

•The new arrival of class $j$ will immediately preempt the lower priority job currently being served and will start its own service.

•When service to the previously preempted class $k$ job eventually resumes (possibly after service to the preempting job of class $j$ and other jobs of priority higher than $k$), *the service will start afresh without remembering the service that has already been provided.*

> *On-going service interrupted by arrival of higher priority.*
>
> *Work already done for the preempted job is not remembered*
>
> ***Work is not conserved***

6

• Arrival Process for Class $i$ is Poisson with rate $\lambda_i$  $i=1,....,P$

• Arrival Processes of different classes independent of each other

• The overall arrival process will also be Poisson with rate $\lambda$

$$\lambda = \sum_{i=1}^{P} \lambda_i$$

7

Service time for Class $i$ has mean $\overline{X}_i$ and second moment $\overline{X_i^2}$ with pdf $b_i(t)$, cdf $B_i(t)$ and L.T. of the pdf as $L_{bi}(s)$

Service times for the different classes assumed to be independent of each other

Traffic of priority class $i$     $\rho_i = \lambda_i \overline{X}_i$       $i=1,.....,P$

Total Traffic  $\rho = \sum_{i=1}^{P} \rho_i = \lambda \overline{X}$

where   $\overline{X} = \sum_{i=1}^{P} \frac{\lambda_i}{\lambda} \overline{X}_i$   is the *mean overall service time*

8

4

Condition for the P-Priority M/G/1 Queue to be Stable

$$r = \sum_{i=1}^{P} r_i < 1$$

For Work-Conserving Queueing Disciplines

For multi-priority queues, it is possible for the queue to become unstable for lower priority traffic while still being stable for the higher priorities.

9

---

**Analytical Approach for Studying Multi-Priority M/G/1 Queues**

Residual Life Approach

*Discussed here and in Sec. 4.5.1 - 4.5.2*

Imbedded Markov Chain

*Discussed in Section 4.5.3*

Also see the additional notes for analytical approaches for Finite Capacity, Multi-Server, Multi-Priority M/M/n/K type Queues

10

**Residual Life Analysis for a Non-Preemptive Priority M/G/1 Queue**

Number of Priority Classes = $P$   (Class 1 lowest priority)

$N_{qk}$   Number of class $k$ jobs waiting in queue (prior to service)

$W_{qk}$   Mean waiting time in queue for jobs of priority class $k$

$N_{qk} = \lambda_k W_{qk}$   (Little's Result for class $k$ jobs)

$R$   Mean Residual Service Time for job currently being served when an arrival (of any priority class) occurs

$$R = \frac{1}{2} \sum_{i=1}^{P} \lambda_i \overline{X_i^2} \qquad (4.40)$$

We now consider each priority class separately, starting with the highest priority class P and ending with the lowest priority class 1

**Class P**

$$W_{qP} = R + \overline{X}_P N_{qP}$$

leading to

$$W_{qP} = \frac{R}{1 - \rho_P} \qquad (4.41)$$

**Class P-1**

$$W_{q(P-1)} = R + \overline{X}_P N_{qP} + \overline{X}_{P-1} N_{q(P-1)} + \overline{X}_P l_P W_{q(P-1)}$$

leading to

$$W_{q(P-1)} = \frac{R}{(1 - r_P)(1 - r_P - r_{P-1})} \qquad (4.44)$$

13

---

**Class P-2**

$$W_{q(P-1)} = R + \overline{X}_P N_{qP} + \overline{X}_{P-1} N_{q(P-1)} + \overline{X}_{P-2} N_{q(P-2)}$$
$$+ \overline{X}_P l_P W_{q(P-2)} + \overline{X}_{P-1} l_{P-1} W_{q(P-2)}$$

leading to

$$W_{q(P-2)} = \frac{R}{(1 - r_P - r_{P-1})(1 - r_P - r_{P-1} - r_{P-2})} \qquad (4.46)$$

14

Therefore, in general, we will get

$$W_{qP} = \frac{R}{1 - r_P} \qquad\qquad i{=}P$$

$$W_{q(P-i)} = \frac{R}{(1 - \sum_{j=0}^{i-1} r_{P-j})(1 - \sum_{j=0}^{i} r_{P-j})} \qquad i{=}1,.....,(P{-}1)$$

(4.47)

$$W_i = W_{qi} + \overline{X}_i \qquad i{=}1,.....,(P{-}1)$$

(4.48)

The parameters $N_i$ and $N_{qi}$ may then be found using Little's Result

15

---

**Residual Life Analysis for a Preemptive Resume Priority M/G/1 Queue**

• Consider $P$ priority classes as before with class $P$ of highest priority

• Jobs of priority classes $1,....., (P{-}1)$ may be interrupted by the arrival of new jobs with higher priority

• No loss of work as interrupted job resumes service from point of interruption

• Queueing Delay can be meaningfully defined only for Class $P$. For the lower priority classes, this parameter will not be important as a job's service can be interrupted even after it starts service

• The Residual Service Time seen by an arrival will depend on the class of the new arrival

16

8

$R_k$ = Mean Residual Service Time as seen by a new job arrival of class $k$

$$R_k = \sum_{i=k}^{P} \frac{1}{2} I_i \overline{X_i^2} \qquad k=1,.....,P \qquad (4.49)$$

• Note that, as mentioned earlier, $R_k$ depends on the class of the new arrival.

• An arrival of the highest priority class will see the smallest mean residual service time as it will preempt any ongoing service of priority class other than itself.

• Arrivals of lower priority class will only be able to preempt jobs of priority lower than themselves

17

**Class P**

*In this case, we can define a mean queueing delay $W_{qP}$ as before*

$$W_{qP} = R_P + \overline{X}_P N_{qP} \qquad \Longrightarrow \qquad W_{qP} = \frac{R_P}{1 - r_P} \qquad (4.50)$$

$$W_P = W_{qP} + \overline{X}_P = \frac{\overline{X}_P(1 - r_P) + R_P}{(1 - r_P)} \qquad (4.51)$$

*Mean Total Delay for Class P*

18

**Class P-1**

$$W_{P-1} = \overline{X}_{P-1} + \frac{R_{P-1}}{1 - r_P - r_{P-1}} + \overline{X}_P l_P W_{P-1} \qquad (4.52)$$

*See Section 4.5.2 for the arguments justifying this term*

This leads to

$$W_{P-1} = \frac{\overline{X}_{P-1}(1 - r_P - r_{P-1}) + R_{P-1}}{(1 - r_P)(1 - r_P - r_{P-1})} \qquad (4.53)$$

*Mean Total Delay for Class P-1*

19

---

In general, we will get

$$W_P = \frac{\overline{X}_P(1 - r_P) + R_P}{(1 - r_P)} \qquad \text{for Class } P$$

$$W_k = \frac{\overline{X}_k(1 - r_P - ...... - r_k) + R_k}{(1 - ...... - r_{k-1})(1 - r_P - ...... - r_k)} \qquad \begin{array}{l}\text{for Class } k, \\ 1 \pounds k \pounds P\text{-}1\end{array}$$

as the total mean delay for each class of customers

20

**Analysis of Multi-Priority M/G/1 Queue using the
Imbedded Markov Chain Approach**

• Though it is possible to do an analysis using this approach
for the work-conserving priority disciplines, this is much
more difficult than the way the mean performance results
were obtained using a Residual Life Approach

• See Section 4.5.3 for the analysis of a *2-Priority M/G/1
Queue* following this approach.

21

11