

# Chapter 4

## M/M/c Queue

### 4.1 Introduction

In Sect. 4.2 below we prove a useful general result (which we call Theorem B) about SP transitions. This theorem facilitates the analysis of *transient* distributions of state variables, and will be applied variously in the sequel.

Sections 4.3–4.5 explain the sample-path structure and dynamics of the generalized M/M/c queue. In the generalized model the SP can make a *transition* between disjoint state-space *sets* (called *pages* or *sheets*). Geometrically, sheets are analogous to a package of sheets of paper, cards in a deck, or pages of a book. They form a discrete number of disjoint subsets of the state space, not connected by a continuous segment of the sample path. (We can also model complex *single-server* queues using the *method of sheets* (aka *method of pages*) (see, e.g., [39, 53, 93]). The *method of sheets* provides LC with great flexibility to analyze different types of stochastic models: e.g., queues; dams (see Sect. 11.8); inventories; production-inventories; actuarial risk models, replacement models; models in the natural sciences, etc.

Sections 4.6.1–4.6.9 develop equations for transient and steady-state pdfs of wait in the generalized M/M/c model. Sections 4.7–4.12 provide steady-state analyses of M/M/c variants using LC. In particular, Sect. 4.8 derives known results for the standard M/M/c queue as a special case of the generalized model. The remaining Sections of the Chapter study variants of M/M/c queues. All Sections provide empirical background for potentially novel applications of LC.

## 4.2 Theorem B for Transient Analysis

We state and prove Theorem B. This straightforward theorem facilitates the transient analyses of a variety of stochastic models.

### 4.2.1 Theorem B

We first give a fundamental generalization of Theorems 3.1 and 3.2 of Chap. 3, which is useful for LC derivations of integro-differential equations for transient distributions in general.

Let  $\{X(t)\}_{t \geq 0}$  denote a sample path of a general stochastic process with state space  $\mathcal{S}$ . Let  $\mathbf{A}, \mathbf{B}$  be arbitrary measurable subsets of  $\mathcal{S}$ . Denote the transient probability  $P(X(t) \in \mathbf{A})$  at instant  $t$  by  $P_t(\mathbf{A})$ ,  $t \geq 0$ . Let  $P_{t_1, t_2}(\mathbf{A}, \mathbf{B})$  be the joint probability  $P(X(t_1) \in \mathbf{A}, X(t_2) \in \mathbf{B})$  at instants  $t_1, t_2 \geq 0$ . Let  $\mathcal{I}_t(\mathbf{A})$  be the number of SP *entrances* and  $\mathcal{O}_t(\mathbf{A})$  the number of SP *exits* of  $\mathbf{A}$ , during  $(0, t)$  (see Fig. 2.7). Assume the derivatives

$$\frac{\partial}{\partial t} P_t(\mathbf{A}), \quad \frac{\partial}{\partial t} E(\mathcal{I}_t(\mathbf{A})), \quad \frac{\partial}{\partial t} E(\mathcal{O}_t(\mathbf{A}))$$

exist for all  $t > 0$ . Both

$$\frac{\partial}{\partial t} E(\mathcal{I}_t(\mathbf{A})) > 0 \quad \text{and} \quad \frac{\partial}{\partial t} E(\mathcal{O}_t(\mathbf{A})) > 0$$

hold wherever the derivatives exist, since  $\mathcal{I}_t(\mathbf{A})$  and  $\mathcal{O}_t(\mathbf{A})$  are counting processes which increase (wide sense, i.e., not strictly; they may be step functions) as  $t$  increases. The following useful result holds.

**Theorem 4.1** *Theorem B* (P.H. Brill, 1983)

$$E(\mathcal{I}_t(\mathbf{A})) = E(\mathcal{O}_t(\mathbf{A})) + P_t(\mathbf{A}) - P_0(\mathbf{A}) \quad (4.1)$$

$$\frac{\partial}{\partial t} E(\mathcal{I}_t(\mathbf{A})) = \frac{\partial}{\partial t} E(\mathcal{O}_t(\mathbf{A})) + \frac{\partial}{\partial t} P_t(\mathbf{A}). \quad (4.2)$$

**Proof** We give two proofs in order to develop intuition about the result.

**Proof 1:** This proof is similar to that of Theorems 3.1 and 3.2 in Sect. 3.2.3 above. We make the correspondence:

$$\begin{aligned} \mathbf{A} &\leftrightarrow (-\infty, x], & \mathcal{I}_t(\mathbf{A}) &\leftrightarrow \mathcal{D}_t(x), & \mathcal{O}_t(\mathbf{A}) &\leftrightarrow \mathcal{U}_t(x), \\ P_t(\mathbf{A}) &\leftrightarrow F_t(x), t \geq 0, & P_{t_1, t_2}(\mathbf{A}, \mathbf{A}) &\leftrightarrow F_{t_1, t_2}(x, x), t_1, t_2 \geq 0. \end{aligned}$$

SP down- and upcrossings of level  $x$  are entrances and exits of sets (Definitions 2.2–2.5). Note that

$$\begin{aligned} \mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A}) = +1 &\iff X(0) \in \mathbf{A}^c, X(t) \in \mathbf{A}, \\ \mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A}) = -1 &\iff X(0) \in \mathbf{A}, X(t) \in \mathbf{A}^c, \\ \mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A}) = 0 &\iff X(0) \in \mathbf{A}, X(t) \in \mathbf{A} \\ &\text{or } X(0) \in \mathbf{A}^c, X(t) \in \mathbf{A}^c. \end{aligned}$$

We thus obtain the following values and corresponding probabilities:

$\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A})$	Probability
+1	$P_{0,t}(\mathbf{A}^c, \mathbf{A}) = P_t(\mathbf{A}) - P_{0,t}(\mathbf{A}, \mathbf{A})$
-1	$P_{0,t}(\mathbf{A}, \mathbf{A}^c) = P_0(\mathbf{A}) - P_{0,t}(\mathbf{A}, \mathbf{A})$
0	$1 - P_t(\mathbf{A}) - P_0(\mathbf{A}) + 2P_{0,t}(\mathbf{A}, \mathbf{A})$

Taking the expected value  $E(\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A}))$  and then the derivative  $\frac{\partial}{\partial t} E(\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A}))$  yields (4.1) and (4.2).

**Proof 2:** Fix  $t \geq 0$ . The probability of the sure event  $\mathbf{S}$  is

$$P_t(\mathbf{S}) = P_t(\mathbf{A} \cup \mathbf{A}^c) = P_t(\mathbf{A}) + P_t(\mathbf{A}^c) = 1.$$

Consider  $P_{t_1, t_2}(\mathbf{A}, \mathbf{S})$ . Events  $\{X(t_1) \in \mathbf{A}\}$  and  $\{X(t_2) \in \mathbf{S}\}$  are independent for every  $0 \leq t_1 \neq t_2$ . Knowledge that  $\{X(t_1) \in \mathbf{A}\}$  has occurred, does not effect the probability of event  $\{X(t_2) \in \mathbf{S}\}$ , which is  $P_{t_2}(\mathbf{S}) = 1$ , and vice versa. Similarly, the events  $\{X(t_1) \in \mathbf{S}\}$  and  $\{X(t_2) \in \mathbf{B}\}$  are independent. Note that  $\mathbf{S} = \mathbf{A} \cup \mathbf{A}^c = \mathbf{B} \cup \mathbf{B}^c$ . Hence

$$\left. \begin{aligned} P_{t_1}(\mathbf{A}) &= P_{t_1, t_2}(\mathbf{A}, \mathbf{S}) = P_{t_1, t_2}(\mathbf{A}, \mathbf{B} \cup \mathbf{B}^c), \\ P_{t_2}(\mathbf{B}) &= P_{t_1, t_2}(\mathbf{S}, \mathbf{B}) = P_{t_1, t_2}(\mathbf{A} \cup \mathbf{A}^c, \mathbf{B}), \end{aligned} \right\}$$

or

$$\left. \begin{aligned} P_{t_1}(\mathbf{A}) &= P_{t_1, t_2}(\mathbf{A}, \mathbf{B}) + P_{t_1, t_2}(\mathbf{A}, \mathbf{B}^c), \\ P_{t_2}(\mathbf{B}) &= P_{t_1, t_2}(\mathbf{A}, \mathbf{B}) + P_{t_1, t_2}(\mathbf{A}^c, \mathbf{B}). \end{aligned} \right\} \tag{4.3}$$

The possible values of  $\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A})$  and corresponding joint probabilities at time points  $t_1 = 0$  and  $t_2 = t > 0$  are:

$\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A})$	Probability
0	$P_{0,t}(\mathbf{A}, \mathbf{A}) + P_{0,t}(\mathbf{A}^c, \mathbf{A}^c)$
+1	$P_{0,t}(\mathbf{A}^c, \mathbf{A})$
-1	$P_{0,t}(\mathbf{A}, \mathbf{A}^c)$

(4.4)

Taking the expected value of  $\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A})$  in (4.4) yields

$$\begin{aligned} E(\mathcal{I}_t(\mathbf{A}) - \mathcal{O}_t(\mathbf{A})) &= P_{0,t}(A^c, \mathbf{A}) - P_{0,t}(\mathbf{A}, A^c) \\ &= P_{0,t}(A^c, \mathbf{A}) + P_{0,t}(\mathbf{A}, \mathbf{A}) \\ &\quad - (P_{0,t}(\mathbf{A}, \mathbf{A}) + P_{0,t}(\mathbf{A}, A^c)) \\ &= P_t(\mathbf{A}) - P_0(\mathbf{A}), \end{aligned}$$

which gives (4.1). Taking  $\partial/\partial t$  in (4.1) yields (4.2). ■

**Remark 4.1** Theorem B also applies to multi-dimensional processes with state space  $\mathcal{S} \subseteq \mathbb{R}^n$ ,  $n = 2, \dots$ , whose states are described by more than one continuous random variable. (**Note:** The symbol  $\mathbb{R}$  denotes the set of real numbers.) We analyze two multi-dimensional inventory models in steady state, in Chap. 7.

### 4.3 Generalized M/M/c Queue

Customers arrive at an M/M/c queue in a Poisson stream at rate  $\lambda$ . There is one waiting line and  $c$  servers. Arrivals start service from the first available server, in order of arrival. We assume that for each arrival, the service time is exponentially distributed with rate selected from a nonempty set  $\boldsymbol{\mu} = \{\mu_0, \dots, \mu_J\}$  of  $J + 1$  positive constants, depending on a *server-assignment policy* specified for the model; this allows service rates to be state dependent. Thus the standard M/M/c queue is a special case (see Sect. 4.8 below; p. 66ff in [84]).

For the generalized M/M/c queue we use a ‘partition/synthesis’ technique. We partition the state space into zero-wait and positive-wait states, and analyze the partitioned states to obtain ‘partial’ pdfs of the waiting time using LC. Then we synthesize those results to obtain the ‘total’ pdf of the waiting time, and related quantities of interest.

We next discuss: virtual wait; server workload; system configuration; the system point process (SP process); and give examples. (References for this section are [11] and [52], and others cited below.)

#### 4.3.1 Virtual Wait and Server Workload

**Notation 4.2** In the remainder of Sect. 4.3 we use two symbols for customers arriving to the system, depending on the context. (1)  $C(t)$  denotes a **would-**

**be** (potential) time- $t$  arrival,  $t \geq 0$ . (2)  $C_{a,t}$  denotes an **actual** time- $t$  arrival, that arrives at  $t^-$ .

Let  $C(t)$  be a would-be (potential) time- $t$  arrival to the system,  $t \geq 0$ . Let  $R_i(t)$  denote the (remaining) *workload* (in time units) at instant  $t \geq 0$ , of server  $i$ ,  $i = 1, \dots, c$  (server numbering is arbitrary). Let  $\{W(t)\}_{t \geq 0}$  be the *virtual wait process*. The random variable  $W(t)$  is the would-be wait required by  $C(t)$  measured from time  $t$  until the start of service of  $C(t)$ . Thus  $W(t) = \min_{i=1, \dots, c} \{R_i(t)\}$ ,  $t \geq 0$ . We assume: sample paths of  $\{W(t)\}_{t \geq 0}$  and of  $\{R_i(t)\}_{t \geq 0}$  are right continuous and have left limits; the model parameters are such that the steady state exists (condition relaxed for the transient analysis in Sects. 4.6.1–4.6.8 below).

**Remark 4.2** In M/M/c ( $c \geq 2$ ), virtual wait  $\neq$  system workload. The system workload at time  $t$  is  $\sum_{i=1}^c R_i(t)$ .

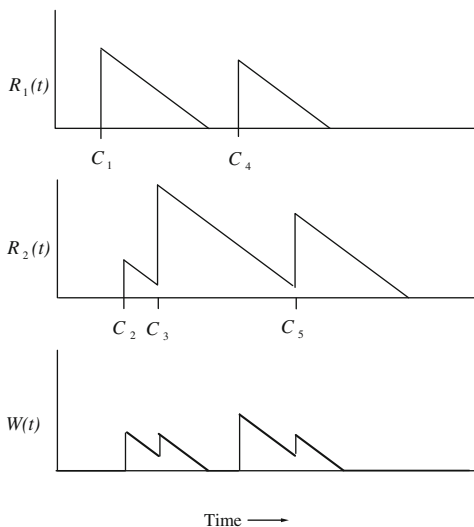
Since  $C_{a,t}$  is an *actual* time- $t$  arrival to the system, that arrives at  $t^-$ ,  $C_{a,t}$ 's wait is  $W(t^-)$  before starting service, from some server  $i_t^*$ . If  $R_j(t^-) = 0$  for some server  $j$ , then  $W(t^-) = 0$ , and  $i_t^*$  is one of the idle servers at  $t^-$ . For zero-wait arrivals, server  $i_t^*$  is selected from the idle servers according to the model's server-assignment policy (e.g., randomly, or by server number, etc.). If  $R_i(t^-) > 0$ ,  $i = 1, \dots, c$ , then  $W(t^-) > 0$ , and  $C_{a,t}$  will start service from server  $i_t^*$  at instant  $t + W(t^-)$  if  $i_t^*$  has the minimum workload among the  $c$  occupied servers at  $t^-$  (further explained in Sect. 4.4 below).

### 4.3.2 Sample Paths of Workload and Virtual Wait

In some models, sample paths of  $R_i(t)$ ,  $i = 1, \dots, c$ , are useful for the overall analysis. We now outline how to construct a sample path of each  $R_i(t)$ ,  $t \geq 0$ ,  $i = 1, \dots, c$ . (Refer to Fig. 4.1, which depicts sample paths for a *special case* of generalized M/M/c with  $c = 2$ .) Without loss of generality, assume the system is empty at  $t = 0$ . Then  $R_i(t) = 0$ ,  $i = 1, \dots, c$ , from  $t = 0$  until the first arrival instant ( $C_1$ ). A new arrival ( $C_2$ ) starts service from a server  $i^*$  and  $R_{i^*}(\cdot)$  jumps upward to the ordinate  $\stackrel{dis}{=} \text{Exp}_{\mu_{i^*}}$ , where  $\mu_{i^*} \in \boldsymbol{\mu}$ .  $R_{i^*}(\cdot)$ , then decreases steadily with slope  $= -1$  as service progresses.

Eventually all  $c$  servers become occupied simultaneously (just after  $C_2$  arrives). Let  $t_1 := \min\{t \mid \text{all } c \text{ servers are occupied}\}$ . If the next customer  $C_\tau$  arrives at time  $\tau > t_1$  before any further service completion, then  $C_\tau$  is the sole customer waiting to start service at time  $\tau$  ( $C_3$ ).  $C_\tau$ 's server will be  $i_\tau^*$  if  $R_{i_\tau^*}(\tau^-) = \min_{i=1, \dots, c} \{R_i(\tau^-)\} := W(\tau^-)$  (virtual wait). The workload

**Fig. 4.1** Sample paths of  $R_j(t)$ ,  $j = 1, 2$ , and  $W(t)$ ,  $t \geq 0$  in M/M/2.  $C_n$ ,  $n = 1, 2, \dots$ , denote customers at successive arrival instants



$R_{i_\tau^*}^*(\tau^-)$  jumps upward by  $C_\tau$ 's service time  $s_{i_\tau^*} = \text{Exp}_{dis}^{\mu_{i_\tau^*}}$ , where  $\mu_{i_\tau^*} \in \boldsymbol{\mu}$ . Thus  $R_{i_\tau^*}^*(\tau) = R_{i_\tau^*}^*(\tau^-) + s_{i_\tau^*} = W(\tau^-) + s_{i_\tau^*}$ . For all other servers,  $R_i(\tau) = R_i(\tau^-)$ ,  $i \neq i_\tau^*$ . Subsequently  $W(\tau) = \min_{i=1, \dots, c} \{R_i(\tau)\}$ .

The next arrival that “sees” at least one idle server ( $C_4$ ), will cause the  $\{W(t)\}_{t \geq 0}$  to evolve similarly. The next arrival that finds all servers busy will be assigned to that server which has minimum workload ( $C_5$ ) and so forth. If arrivals find several customers waiting in line, the dynamics are similar to the case ‘all servers busy’ (described in Sect. 4.4 below).

### 4.3.3 Distinguishable Servers

When tracking server workloads, we regard the servers as distinguishable (Fig. 4.1). However, we are often interested in the statistical properties of the entire system, rather than the processing of each individual customer, or the action of a particular server. Here we analyze the system by constructing a sample path of  $\{W(t)\}_{t \geq 0}$  generated according to the model’s prescribed probability laws for the service and interarrival times, and operational policies.

Suppose we can keep track of the  $c$  server workloads in continuous time. Then we could assign a ‘ticket’ to each arrival, which points to its up-coming server, identified because it has the minimum workload at the arrival instant.

This procedure distributes *theoretical* waiting lines to the  $c$  servers, although there is only one physical waiting line in the waiting room.

#### 4.3.4 Indistinguishable Servers

In many M/M/c models, it is not necessary to construct sample paths of the server workloads  $\{R_i(t)\}_{t \geq 0}$ ,  $i = 1, \dots, c$ , in order to construct a sample path of  $\{W(t)\}_{t \geq 0}$ . It suffices to regard the servers as *indistinguishable*. Then it is not necessary to track individual server workloads. To analyze important statistical properties of the model, it is sufficient to track *directly* the virtual wait  $W(t) := \min_{i=1, \dots, c} \{R_i(t)\}$ . Thus we utilize what we call the *system configuration* (Sect. 4.4).

### 4.4 System Configuration

In generalized M/M/c, assume a ‘system manager’ knows the up-coming target server  $i_t^*$  to be occupied at instant  $t + W(t^-)$  by a would-be time- $t$  arrival, denoted by  $C(t)$  (i.e., the manager knows the server having minimum workload at time  $t$ ). Let  $\mathbf{M}(t) :=$  *system configuration* at time  $t$ . The process  $\{\mathbf{M}(t)\}_{t \geq 0}$  tracks the service rates of the  $c - 1$  servers *other than*  $i_t^*$ . We assume that the model specifies  $J + 1$  possible exponential service rates:  $\boldsymbol{\mu} = \{\mu_0, \mu_1, \dots, \mu_J\}$ . Each arrival is assigned a service rate selected from the set  $\boldsymbol{\mu}$ . Recall that if  $t$  is not an arrival instant, sample-path right continuity implies  $W(t^-) = W(t)$ .

**Definition 4.1** The system configuration  $\mathbf{M}(t)$  is a  $J + 1$  vector of **server occupancy numbers**  $m_j \geq 0$ , namely  $\mathbf{M}(t) = (m_0, \dots, m_J)$ , where  $m_j :=$  **number of servers having service rate**  $\mu_j \in \boldsymbol{\mu}$ , *among the  $c - 1$  servers other than  $i_t^*$  at  $t + W(t^-)$ .*

**Definition 4.2** The set of all possible configurations is denoted by  $\mathbf{M} = \{\mathbf{m} | \mathbf{m} = (m_0, \dots, m_J)\}$ .

For each configuration  $\mathbf{m} \in \mathbf{M}$ ,  $0 \leq \sum_{j=0}^J m_j \leq c - 1$ ;  $C(t)$  would start service at instant  $t + W(t^-)$  and would be assigned a service rate  $\mu_t(W(t^-), \mathbf{m}) \in \boldsymbol{\mu}$ , which may be a function of three variables:  $t$ ,  $W(t^-)$ , and  $\mathbf{m}$ . That is

$$(t, W(t^-), \mathbf{m}) \rightarrow \mu_j \in \boldsymbol{\mu}, \text{ for some } j = 0, \dots, J.$$

**Remark 4.3** In various models, the service rate  $\mu_t(\cdot, \cdot)$  may also depend on other variables as well. It may be selected randomly from the set  $\mu$ . Additionally, the number of possible service rates in  $\mu$  may be countable.

#### 4.4.1 Inter Start-of-service Depart Time $\mathcal{S}_t$

In generalized M/M/c a *key random variable* is the time- $t$  ‘look-ahead’ **inter start-of-service depart time**, denoted by  $\mathcal{S}_t$ . For example, let the state  $(W(t^-), M(t^-))$  be  $(x, \mathbf{m})$  when customer  $C_{a,t}$  arrives. Then  $C_{a,t}$ ’s required wait before starting service is  $x \geq 0$  and the configuration is  $\mathbf{m}$ . If  $0 \leq \sum_{i=0}^J m_j \leq c - 1$  then  $x = 0$  and  $C_{a,t}$  starts service immediately by one of the idle servers. If  $\sum_{i=0}^J m_j = c - 1$  there are two possibilities for  $x$ . If  $x = 0$  then  $C_{a,t}$  starts service immediately by the *single idle server*. If  $x > 0$  then  $C_{a,t}$  waits time  $x$  before starting service by the first available server thereafter. Just after  $C_{a,t}$  starts service at time  $t + x$  all  $c$  servers will be occupied.

**Definition 4.3** The *inter start-of-service depart time*  $\mathcal{S}_t$  is the time measured from  $t + x$  (start of  $C_{a,t}$ ’s service time) until the first departure from the system after  $t + x$ . In other words  $\mathcal{S}_t :=$  *time from the start of service of  $C_{a,t}$  until the first departure from the system thereafter*.

Importantly,

$$\mathcal{S}_t = \min_{dis} \{ \text{Exp}_{m_0 \mu_0}, \dots, \text{Exp}_{m_J \mu_J}, \text{Exp}_{\mu_t(W(t^-), \mathbf{m})} \},$$

which is the *minimum* of  $\sum_{i=0}^J m_j + 1 (= c)$  independent exponential r.v.s. Among these,  $m_j$  servers have rate  $\mu_j$ ,  $j = 0, \dots, J$ , and one server has rate  $\mu_t(x, \mathbf{m})$  (assigned to  $C_{a,t}$ ). Thus  $\mathcal{S}_t = \text{Exp}_{\nu_t}$  where  $\nu_t = \sum_{j=0}^J m_j \mu_j + \mu_t(x, \mathbf{m})$ .

An important aspect of the forgoing definition of system configuration and use of  $\mathcal{S}_t$  when  $W(t^-) > 0$ , is that once all  $c$  servers are occupied, the probability distribution of  $\mathcal{S}_t$  is independent of future arrivals to the system. That is, the set of active servers functions like a separate sub-system until the first departure thereafter, mindful of the memoryless property of the exponential service times. Although the concept ‘*system configuration*’ may appear ‘different’, it is straightforward to apply when developing model equations for the pdf of the waiting time in complex M/M/c queues (see, e.g., pp. 80–97 in [11], and also in [38, 52, 53]).



### 4.4.2 Number of Configurations

Let  $(W(t), M(t)) = (x, \mathbf{m})$  (see Definition 4.1 above). Looking ahead to  $t + W(t)$ , assume  $\mathbf{m} = (m_0, \dots, m_J)$  is such that

$$\sum_{j=0}^J m_j = k, 0 \leq k \leq c - 1.$$

The servers are considered to be *indistinguishable* (as in subsequent models in this monograph, unless otherwise noted). We track only the *number of servers occupied* with service rate  $\mu_j \in \boldsymbol{\mu}, j = 0, \dots, k$ .

The number of possible configurations such that exactly  $k$  servers are occupied, is the number of non-negative integer solutions of the equation

$$m_0 + \dots + m_J = k.$$

It is the same as the number of ways of distributing  $k$  *indistinguishable* balls in  $J + 1$  *distinguishable* cells, namely  $\binom{J+k}{J} = \binom{J+k}{k}$  (see Lemma, p. 36, Chap. II in [73]). Thus, the total number of possible configurations is

$$\sum_{k=0}^{c-1} \binom{J+k}{J} = \binom{J+c}{J+1} = \binom{J+c}{c-1}. \tag{4.5}$$

The first equality in (4.5) is readily proved by induction.

**Example 4.1** Consider an M/M/ $c$  queue with  $c = 3$  and  $J = 2$ , so that  $\boldsymbol{\mu} = \{\mu_0, \mu_1, \mu_2\}$ . If a potential arrival  $C(t)$  finds the system **empty**, then  $(W(t), M(t)) = (0, (0, 0, 0))$ ; thus  $m_i = 0, i = 0, 1, 2$ .  $C(t)$  would wait zero and “see” zero servers occupied,  $(0, 0, 0)$ . The number of solutions of  $m_0 + m_1 + m_2 = 0$  is  $\binom{J+0}{J} = \binom{2}{2} = 1$ .  $C(t)$  would wait  $W(t^-) = 0$  and start service from one the three unoccupied servers, per the server-assignment policy.

If  $C(t)$  would find **one** customer in the system (one occupied server), then  $W(t^-) = 0$  and the configuration that  $C(t)$  would “see” is one of **three** possible vectors

$$M(t^-) \in \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}.$$

The number of solutions of  $m_0 + m_1 + m_2 = 1$  is  $\binom{J+1}{J} = \binom{3}{2} = 3$ .  $C(t)$  would start service from one of the two unoccupied servers, per the server-assignment policy.

If  $C(t)$  would find **two** customers in the system, then  $W(t^-) = 0$  and the configuration that  $C(t)$  would “see” is one of **six** possible vectors

$$M(t^-) \in \{(2, 0, 0), (0, 2, 0), (0, 0, 2), (1, 1, 0), (1, 0, 1), (0, 1, 1)\}.$$

The number of solutions of  $m_0 + m_1 + m_2 = 2$  is  $\binom{J+2}{J} = \binom{4}{2} = 6$ .  $C(t)$  would start service from the **one unoccupied** server.

If  $C(t)$  would find **three or more** customers in the system, then all **three** servers would be occupied at  $t^-$ . The look-ahead configuration that  $C(t)$  would “see” **just before start of service** at  $t + W(t^-)$  is also one of **six** possible vectors

$$M(t) \in \{(2, 0, 0), (0, 2, 0), (0, 0, 2), (1, 1, 0), (1, 0, 1), (0, 1, 1)\}.$$

The six possible configurations are the *same* as when  $C(t)$  sees **two occupied** servers. This is because a configuration tracks the service-rate occupancies of those servers **other than**  $C(t)$ 's eventual server. Customer  $C(t)$  would wait a **positive time** and start service at  $t + W(t^-)$  from some server  $i_t^*$ . We “look ahead” to the **start of service instant**  $t + W(t^-)$  **and assign rate**  $\mu_t(W(t^-), M(t^-))$  to  $i_t^*$ . The random variable  $M(t)$  tracks the service-rate occupancies of the two servers other than  $i_t^*$  at  $t + W(t^-)$ . (The look-ahead idea is not new. For example, it is tacitly assumed for the virtual wait in the standard M/G/1 queue, where we increase the virtual wait by a service time at an *arrival* instant, although the service is not started until the end of the waiting time. The generalized M/M/c generalizes the M/G/1 look-ahead idea to the start-of-service time.)

At instant  $t$ , the state  $(W(t), M(t))$  conveys sufficient information to determine the probabilities of the  $m_j$  ( $j = 0, \dots, J$ ) occupied servers that will have the minimum service time among all the occupied servers at time  $t + W(t^-)$ . These probabilities depend on the Markovian property. We shall illustrate this more fully in Example 4.2 below.

The total number of possible configurations is

$$\sum_{k=0}^{c-1} \binom{J+k}{J} = \sum_{k=0}^2 \binom{J+k}{J} = \binom{J+c}{J+1} = \binom{5}{3} = 10.$$

### 4.4.3 Border States

In Example 4.1 the zero-wait state  $\{(0, \mathbf{m})\}$  is a **border** state if

$$\mathbf{m} \in \{(2, 0, 0), (0, 2, 0), (0, 0, 2), (1, 1, 0), (1, 0, 1), (0, 1, 1)\}.$$

**Definition 4.4** We call the state  $(W(t), M(t))$  a **border state** if  $W(t) = 0$  and  $M(t)$  is such that  $\sum_{j=0}^J m_j = c - 1$ . A border state is a discrete zero-wait state in a **boundary** separating other discrete zero-wait states and a set of continuous positive-wait states.

In the above definition, the *other* zero-wait states are *non-border* states such that  $0 \leq \sum_{j=0}^J m_j < c - 1$ . Border states communicate with continuous positive-wait states in one step: at arrival instants (zero-wait  $\rightarrow$  positive-wait); or at departure instants (positive-wait  $\rightarrow$  zero-wait). When the SP moves on a path *from* a *non-border* zero-wait state *to* a continuous positive-wait state the path must pass through a border state at an arrival instant. In the opposite direction, *from* a positive-wait state *to* a non-border zero-wait state, the path must pass through a border state at a departure instant.

We denote the set of border states by  $S_b$ , and the set of border configurations by  $M_b$ . Thus

$$\begin{aligned} S_b &= \left\{ (0, \mathbf{m}) \mid \sum_{j=0}^J m_j = c - 1 \right\}, \\ M_b &= \{ \mathbf{m} \mid (0, \mathbf{m}) \in S_b \} = \left\{ \mathbf{m} \mid \sum_{j=0}^J m_j = c - 1 \right\}. \end{aligned} \quad (4.6)$$

#### 4.4.4 The Next Configuration

Consider an actual arrival  $C_{a,t}$  at instant  $t$ .  $C_{a,t}$  “sees” configuration  $M(t^-)$ . Just *after* the arrival the configuration is  $M(t)$ . Either  $M(t) = M(t^-)$  or  $M(t) \neq M(t^-)$ . We illustrate by example how to compute the probability mass function of  $M(t)$ .

**Example 4.2** Consider Example 4.1 for M/M/c with  $c = 3$ ,  $J = 2$ . Suppose  $C_{a,t}$  arrives when the wait is  $W(t^-)$  and the configuration is  $(m_0, m_1, m_2) = (1, 1, 0)$ . The state is

$$(W(t^-), M(t^-)) = (0, (1, 1, 0)).$$

Suppose that  $C_{a,t}$  is assigned service rate  $\mu_0$ , i.e.,  $\mu_t(W(t^-), (1, 1, 0)) = \mu_0$ .

At instant  $t + 0$ , **just after**  $C_{a,t}$  **starts service**, there will be **two** servers with rate  $\mu_0$  since  $m_0 = 1$ . There will be *one* server with rate  $\mu_1$ , since  $m_1 = 1$ . The inter-start-of-service-depart time  $S_t = \text{Exp}_{dis}^{2\mu_0 + \mu_1}$ .

We now compute the probability distribution of the **next configuration** at instant  $t + W(t)$ . Thus,

$$\begin{aligned}
 P(M(t) = (2, 0, 0)) & \\
 &= P(\text{rate-}\mu_1 \text{ server finishes first}) \\
 &= \frac{\mu_1}{2\mu_0 + \mu_1},
 \end{aligned}$$

$$\begin{aligned}
 P(M(t) = (1, 1, 0)) & \\
 &= P(\text{rate-}\mu_0 \text{ server finishes first}) \\
 &= \frac{2\mu_0}{2\mu_0 + \mu_1}.
 \end{aligned}$$

Importantly

$$P(M(t) = (2, 0, 0)) + P(M(t) = (1, 1, 0)) = 1.$$

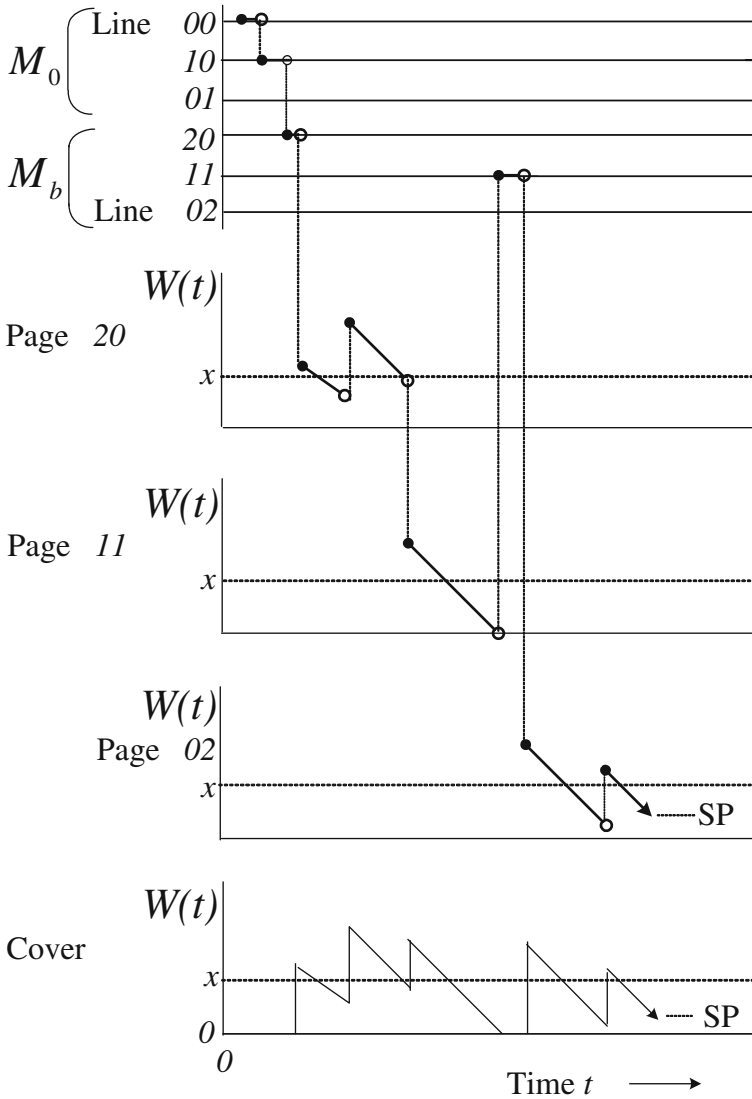
The only two possible configurations for  $M(t)$  are  $(2, 0, 0)$  and  $(1, 1, 0)$ , independent of whether  $W(t^-) = 0$  or  $W(t^-) > 0$  (illustrated below in Example). No other configuration is possible for  $M(t)$  once the arrival at  $t^-$  has been assigned rate  $\mu_0$ . Knowledge of their probabilities is sufficient to analyze the sample path to write down model equations using LC.

**Remark 4.4** The service mechanism can be generalized considerably. We can expand the domain of  $\mu_t(w, \mathbf{m})$  to include: type or priority class of  $C_{a,t}$ ; type of customer replaced by  $C_{a,t}$  in server  $i_t^*$ ; type of any customer followed by  $C_{a,t}$  into service; identity of server  $i_t^*$  (e.g., server number or unique property); number of customers in the system or waiting for service at the arrival or start of service instant of  $C_{a,t}$ ; various types of bounds on the virtual wait; reneging indices; blocked and cleared customers, etc. (see, e.g., [38] for an effective definition of  $M(t)$  due to L. Green; [42]; also see [39, 44, 53]; and others).

Other generalizations may incorporate: a non-homogeneous Poisson arrival process with intensity  $\lambda_t$ , or a Poisson arrival rate  $\lambda(W(t), M(t))$  which is a function of the current state  $(W(t), M(t))$ ; or various Markov arrival processes.

## 4.5 System Point Process

We now discuss the *system point process* and the geometry of its state space (see Fig. 4.2).



**Fig. 4.2** Sample path of SP process  $\{W(t), M(t)\}_{t \geq 0}$  for Example 4.4 with random assignment of service rates independent of state at arrival instants ( $c = 3, J = 1$ ). The space  $T \times S$  has 6 lines for zero-wait states, and 3 pages (sheets) for positive-wait states (pages 20, 11, 02). The cover is the projection of the sample path from all lines and pages onto one non-negative planar quadrant

We call  $\{W(t), M(t)\}_{t \geq 0}$  the *system point (SP) process*. Its nomenclature derives from the fact that the SP traces out a sample path as the system evolves over time. The SP process for M/M/c queues is a generalization (with

exponential service times) of the virtual wait process for M/G/1 queues. State variable  $W(t) := \text{virtual wait}$ ; state variable  $M(t) := \text{system configuration at time } t$ . Random variable  $M(t)$  is discrete. The SP process is a Markov process (Sect. 4.5.7 below).

We partition the *state space*  $\mathcal{S}$  into three disjoint state-space sets  $\mathcal{S}_0$ ,  $\mathcal{S}_b$ ,  $\mathcal{S}_1$ .  $\mathcal{S}_0$  contains the zero-wait states that are non-boundary states.  $\mathcal{S}_b$  contains the zero-wait states that are boundary states.  $\mathcal{S}_1$  contains the positive-wait states. The states in  $\mathcal{S}_0 \cup \mathcal{S}_b$  are atoms (see Sect. 2.4.9 for ‘atom’). The states in  $\mathcal{S}_1$  are points in a continuum, e.g.,  $(x, \mathbf{m})$ ,  $x > 0$ , and  $\mathbf{m} = (m_0, \dots, m_J)$ . Specifically,

$$\begin{aligned}\mathcal{S}_0 &= \{(0, \mathbf{m}) \mid 0 \leq \sum_{j=0}^J m_j \leq c - 2\}, \\ \mathcal{S}_b &= \{(0, \mathbf{m}) \mid \sum_{j=0}^J m_j = c - 1\}, \\ \mathcal{S}_1 &= \{(x, \mathbf{m}) \mid x > 0, \sum_{j=0}^J m_j = c - 1\}.\end{aligned}\tag{4.7}$$

Note that  $\mathcal{S} = \mathcal{S}_0 \cup \mathcal{S}_b \cup \mathcal{S}_1$ , and  $\mathcal{S}_0 \cap \mathcal{S}_b = \mathcal{S}_0 \cap \mathcal{S}_1 = \mathcal{S}_b \cap \mathcal{S}_1 = \phi$ , the empty set. The corresponding sets of system configurations are

$$\begin{aligned}\mathbf{M}_0 &= \{\mathbf{m} \mid (0, \mathbf{m}) \in \mathcal{S}_0\} = \{\mathbf{m} \mid 0 \leq \sum_{j=0}^J m_j \leq c - 2\}; \\ \mathbf{M}_b &= \{\mathbf{m} \mid (0, \mathbf{m}) \in \mathcal{S}_b\} = \{\mathbf{m} \mid \sum_{j=0}^J m_j = c - 1\}; \\ \mathbf{M}_1 &= \{\mathbf{m} \mid (x, \mathbf{m}) \in \mathcal{S}_1\} = \{\mathbf{m} \mid x > 0, \sum_{j=0}^J m_j = c - 1\};\end{aligned}\tag{4.8}$$

thus  $\mathbf{M}_b = \mathbf{M}_1$  (see Sect. 4.4 above).

**$W(t^-) = 0$**  An arrival  $C(t)$  would “see”  $W(t^-) = 0$  if and only if the state at time  $t^-$  is in  $\mathcal{S}_0 \cup \mathcal{S}_b$ .  $C(t)$  would then wait zero and start service from some server, say  $i_t^*$ , at time  $t$ . Geometrically, we associate a distinct horizontal line  $T \times (0, \mathbf{m})$  with each state  $(0, \mathbf{m}) \in \mathcal{S}_0 \cup \mathcal{S}_b$  where  $T$  is the time axis  $[0, \infty)$ . We call the line  $T \times (0, \mathbf{m})$  “*line*  $\mathbf{m}$ ” (e.g., Fig. 4.2).

**$W(t^-) > 0$**   $C(t)$  would “see”  $W(t^-) > 0$  if and only if the state is in  $\mathcal{S}_1$ .  $C(t)$  would wait time  $W(t^-)$  and start service from some server, say  $i_t^*$ , at time  $t + W(t^-)$ . Geometrically, we associate the quadrant of the plane  $T \times (0, \infty)$  with each set of continuous states  $(x, \mathbf{m}) \in \mathcal{S}_1$ ,  $x > 0$ . We call the positive quadrant  $T \times ((0, \infty), \mathbf{m})$  *sheet* (or *page*)  $\mathbf{m}$  (e.g., Fig. 4.2).

**Sample Path Diagram** The LC analyst draws (or visualizes) a plot of  $W(t)$  versus  $t$  on page  $m$  while the system is in the state corresponding to configuration  $m$ . In a diagram, we may place the zero-wait *border* lines (corresponding to the states in  $S_b$ )  $T \times (0, m)$ ,  $(0, m) \in S_b$ , alongside the zero-wait *non-border* lines for states  $(0, m) \in S_0$ ; or else at level 0 of the corresponding sheets for the positive-wait states  $S_1$  having the same configurations  $m$ . There is a one-to-one correspondence between sheets and states in  $S_b$  (Fig. 4.2).

### 4.5.1 Sample Path of SP Process

A sample path of  $\{W(t), M(t)\}_{t \geq 0}$  is a piecewise right-continuous function of  $t$  having left limits. It has a finite number of jumps during finite time intervals (see Sect. 2.2 and Definition 2.1 in Chap. 2). We plot a sample path within a Cartesian product space  $T \times S = T \times (S_0 \cup S_b \cup S_1)$ . The direction of time is from left to right. It is useful to envisage each Cartesian product  $T \times (0, m)$ ,  $(0, m) \in S_0 \cup S_b$  as a *line*; and each quadrant  $T \times ((0, \infty), m)$ ,  $m \in M_1$ , as a *sheet* (or page in a *book*).

#### Description of a Sample Path

Assume that the system starts empty. The SP moves among the zero-wait lines, jumping from line to line at arrival and departure instants, eventually reaching a zero-wait *border* line (often placed at level zero of some sheet  $m \in M_b$ ). Eventually the SP jumps from line  $m \in M_b$ , to a positive level on some sheet ' $k$ ', at an arrival instant. It then moves steadily with slope  $-1$  on sheet  $k$ . It is possible that either  $m = k$ , or  $m \neq k$ , depending on the probabilities governing the motion of the SP. (See Fig. 4.2 and Example 4.4 in Sect. 4.5.5 for a detailed example.) Other clarifying examples are in the author's Ph.D. thesis (Fig. 4.3, p. 79, Chap. 4 in [11]; and in [52]).

At an arrival instant while the SP is on sheet  $k$ , the SP may jump to another sheet, say  $m'$ , and move steadily with slope  $-1$  on sheet  $m'$  for a positive time. Otherwise the SP may jump, and stay on the same sheet  $k$ . On each sheet it moves downward with slope  $-1$ . If the SP hits level 0 from above on page  $k$  before the next arrival, it starts moving immediately on the border line  $k$  (no customers waiting,  $c - 1$  servers occupied).

If the SP is in a state in  $S_b \cup S_1$  having configuration  $m$  at some arrival instant, it makes a jump ending either on page  $m$  or on some page  $k \neq m$ . If  $k \neq m$ , the SP is said to make an  $m \rightarrow k$  *transition*. This may be an upward

jump from a border line  $m$ , or from sheet  $m$  to sheet  $k$ , at an arrival instant. Generally,  $m \rightarrow k$  transitions do not give rise to ‘typical’ level crossings as in M/G/1 models that have exactly one ‘page’. However,  $m \rightarrow m$  transitions from a border line  $m$  or from a point on sheet  $m$  to a higher point on sheet  $m$ , are similar to SP jumps as discussed for models with a single sheet (Sect. 2.3 in Chap. 2).

**Remark 4.5** In some model variants, an  $m \rightarrow k$  transition may be a **parallel jump**. That is, the SP makes a jump from a level  $y$  on page  $m$  to level  $y$  on page  $k \neq m$ , at an arrival instant. For example, in an M/G/1 queue, we may utilize a modified configuration  $M(t) = n$ , where  $n$  is the number of customers waiting for service, and the virtual wait is unchanged at some arrival instants. Such parallel jumps occur in M/G/1 or M/M/c queues with bulk service (see, e.g., [93]).

#### 4.5.2 A Metaphor for Sample Path and SP Motion

The SP motion over the state space is like the motion of the tip of a pen writing out a *single-page* history of the system over time. The writing takes place in a book of transparent pages all the same size. The cover is also transparent. The pen moves from left to right, and never overlays what has been written already. After writing flat or sloped lines on a page for a random amount of time, the pen jumps to a different page, and continues writing. The pen jumps in this manner at random time points from page to page. The next page is selected by a random mechanism depending on where it is presently. The *entire* history up to an instant in time can be seen only by holding all the pages one behind the other, like pages in a book, and viewing the projected history on the *cover*. The projected history on the cover is invariant with respect to shuffling or mixing the pages, which change their relative positions. An analyst that views an arbitrary page in isolation, sees only local segments of the history specific to that page (see Fig. 3.2, p. 49 and Fig. 4.3, p. 79 in [11]).

The global history is like the *total* sample path of the SP process over the state space  $S_0 \cup S_b \cup S_1$ . The local histories on various pages are like sample-path segments due to sojourns on the ‘lines’ and ‘sheets’ of the state space. On the cover, SP motion on all the lines occurs at level 0. That is, when all the lines are projected onto the cover, they are placed at level zero—to form a single “line 0”.



We may think of the overall method as having several steps.

1. **Partition** the **Time-State space** into mutually exclusive and exhaustive lines and sheets.
2. **Analyze** the sample-path segments on the lines and sheets using LC methods.
3. **Project** the sample-path segments from the lines and sheets onto the ‘cover’ of the ‘book’. Analyze the projected path on the cover using LC.
4. **Combine** all the LC results with a normalizing condition. Construct the model equations (usually Volterra integral equations of the second kind with parameter for the pdf’s of interest) and derive probability distributions of the model.

The LC method utilizes statistical properties of the local path segments on the lines and sheets. It also uses statistical properties of the projected path on the cover. It employs the one-step communication properties among the lines and sheets (at successive arrivals and/or departures) to construct a sample path. Basic LC theorems apply to each page  $m \in M$ . Jumps out of, and into lines and sheets, follow rate-conservation laws.

### Equations for PDF of Wait

We use sample-path structure, and transition rates into and out of state-space sets, to construct (by inspection of the sample path) integro-differential and differential equations in a *transient* analysis. Similarly, we construct integral equations and algebraic equations in a *steady-state* analysis. These are equations for the joint pdf and/or cdf of *wait and system configuration*. We can also derive equations for the marginal (total) pdf and cdf of wait, or for the probabilities of the system configurations.

**Remark 4.6** The author originally had the idea for partitioning the state space, visualizing the positive-wait states over time in separate quadrants, and having a ‘system point’ moving on the quadrants over time, from an analogy with Riemann sheets and diagrams of winding numbers in complex variable theory (see, e.g., Sect. 3.4, p. 137ff in [96]). My Ph.D. thesis used the term ‘sheets’. The term ‘pages’ was introduced soon after. I also thought of using the term ‘cards’, analogous to boxes of computer cards for data and programs, which were widely in use in 1974. Then, the state space could be pictured like a box or deck of rectangular cards. Such (‘IBM’) cards had been ubiquitous until personal computers became common in the 1980s.

### 4.5.3 Notation: Probabilities and Distributions

#### Transient Probabilities and Distributions

We denote: the zero-wait probabilities by

$$P_t(0, \mathbf{m}) = P(W(t) = 0, M(t) = \mathbf{m}), (0, \mathbf{m}) \in S_0 \cup S_b;$$

the mixed joint cdf of  $(W(t), M(t))$  by

$$\begin{aligned} F_t(x, \mathbf{m}) &= P(W(t) \leq x, M(t) = \mathbf{m}) \\ &= P_t(0, \mathbf{m}) + P(0 < W(t) \leq x, M(t) = \mathbf{m}) \\ &= P_t(0, \mathbf{m}) + \int_{y=0}^x f_t(y, \mathbf{m}) dy, \quad x \geq 0, t \geq 0, \\ &\quad (0, \mathbf{m}) \in S_0 \cup S_b, (x, \mathbf{m}) \in S_1, \end{aligned}$$

where  $P(0 < W(t) \leq x, M(t) = \mathbf{m}) = P(\phi) = 0$  if  $x = 0$ ; and the mixed joint pdf of  $(W(t), M(t))$  is

$$f_t(x, \mathbf{m}) = \frac{\partial}{\partial x} F_t(x, \mathbf{m}), \quad x > 0, t \geq 0, (x, \mathbf{m}) \in S_1,$$

wherever  $\frac{\partial}{\partial x} F_t(x, \mathbf{m})$  exists.

We assume:

1.  $F_t(x, \mathbf{m})$  and  $f_t(x, \mathbf{m})$  are right continuous in  $x$  for every  $t \geq 0, \mathbf{m} \in M_1$ .
2.  $\frac{\partial}{\partial t} F_t(x, \mathbf{m})$  and  $\frac{\partial}{\partial t} f_t(x, \mathbf{m}), t > 0, x \geq 0$ , exist and are finite for every  $\mathbf{m} \in M_1$ .

Let  $P_0(t) := P(W(t) = 0)$  be the marginal probability of a zero wait at  $t$ . Then

$$\begin{aligned} P_0(t) &= \sum_{(0, \mathbf{m}) \in S_0 \cup S_b} P_t(0, \mathbf{m}) \\ &= \sum_{(0, \mathbf{m}) \in S_0} P_t(0, \mathbf{m}) + \sum_{(0, \mathbf{m}) \in S_b} P_t(0, \mathbf{m}), \quad t \geq 0. \end{aligned}$$

The transient marginal cdf of wait  $P(W(t) \leq x)$  is

$$\begin{aligned}
 F_t(x) &= \sum_{(0, \mathbf{m}) \in \mathcal{S}_0} P_t(0, \mathbf{m}) + \sum_{(0, \mathbf{m}) \in \mathcal{S}_b} F_t(x, \mathbf{m}) \\
 &= \sum_{(0, \mathbf{m}) \in \mathcal{S}_0 \cup \mathcal{S}_b} P_t(0, \mathbf{m}) + P(0 < W(t) \leq x) \\
 &= P_0(t) + P(0 < W(t) \leq x) \\
 &= P_0(t) + \int_{y=0}^x f_t(y) dy, \quad x \geq 0, t \geq 0.
 \end{aligned}$$

Note that  $P_t(0, \mathbf{m}) = F_t(0, \mathbf{m})$  for  $(0, \mathbf{m}) \in \mathcal{S}_b$ .

(Recall the definitions of  $\mathcal{M}_b$  and  $\mathcal{S}_b$  in (4.6), and  $\mathcal{M}_b = \mathcal{M}_1$ , which is the set of system configurations for positive-wait states.)

The transient marginal pdf of  $W(t)$  is

$$f_t(x) = \frac{\partial}{\partial x} F_t(x) = \sum_{\mathbf{m} \in \mathcal{M}_1} f_t(x, \mathbf{m}), \quad x > 0, t \geq 0.$$

A potential (would-be) arrival  $C(t)$  would find the system configuration to be  $\mathbf{m} \in \mathcal{M}_0 \cup \mathcal{M}_b$  with probability  $P_t(0, \mathbf{m})$ .  $C(t)$  would find the configuration to be  $\mathbf{m} \in \mathcal{M}_1$  with probability  $F_t(\infty, \mathbf{m})$ . The normalizing condition for fixed  $t \geq 0$ , is

$$\begin{aligned}
 F_t(\infty) &= \sum_{\mathbf{m} \in \mathcal{M}_0} P_t(0, \mathbf{m}) + \sum_{\mathbf{m} \in \mathcal{M}_1} F_t(\infty, \mathbf{m}) \\
 &= \sum_{\mathbf{m} \in \mathcal{M}_0 \cup \mathcal{M}_b} P_t(0, \mathbf{m}) + \sum_{\mathbf{m} \in \mathcal{M}_1} \int_{y=0}^{\infty} f_t(y, \mathbf{m}) dy \\
 &= \sum_{(0, \mathbf{m}) \in \mathcal{S}_0 \cup \mathcal{S}_b} P_t(0, \mathbf{m}) + \sum_{\mathbf{m} \in \mathcal{M}_1} \int_{y=0}^{\infty} f_t(y, \mathbf{m}) dy = 1.
 \end{aligned}$$

### Steady-State Probabilities and Distributions

We denote the steady-state zero-wait probabilities, pdfs and cdfs of wait by dropping the subscript  $t$  in the immediately foregoing notation for the transient quantities.

### 4.5.4 Configuration Just After an Arrival

Example 4.3 below demonstrates the probability of a system configuration just after an arrival. Assume that an actual customer  $C_{a,t}$  arrives and finds the state to be  $(W(t^-), \mathbf{M}(t^-)) = (x, \mathbf{m})$ . The service rate assigned to  $C_{a,t}$  is  $\mu_t(x, \mathbf{m}) \in \boldsymbol{\mu}$ . Recall that sample paths are right continuous and have left limits.

**Example 4.3** Consider Example 4.2 in Sect. 4.4.4, where  $c = 3$ ,  $J = 2$ . Let each arrival receive a service rate selected **with equal probability** from the set  $\boldsymbol{\mu} := \{\mu_0, \mu_1, \mu_2\}$ . Then

$$P(C_{a,t} \text{ starts service at } t + W(t^-) \text{ with service rate } \mu_i) = \frac{1}{3}, i = 0, 1, 2,$$

independent of  $t$  and  $W(t^-)$ . Assume  $(W(t^-), \mathbf{M}(t^-)) = (x, (2, 0, 0))$ ,  $x > 0$ , just before  $C_{a,t}$  arrives. Then  $C_{a,t}$  will wait a positive time  $x$ . Looking ahead to time  $t + W(t^-)$ , two other occupied servers will have service rates  $\mu_0$  ( $m_0 = 2$ ,  $m_1 = m_2 = 0$ ) when  $C_{a,t}$  starts service at  $t + W(t^-)$ , in the just-vacated server. **Question:** What is the configuration  $\mathbf{M}(t)$  just after  $C_{a,t}$  arrives? It can be either  $(2, 0, 0)$ ,  $(1, 1, 0)$ , or  $(1, 0, 1)$ . The probabilities for  $\mathbf{M}(t)$  are:

$$\begin{aligned} P(\mathbf{M}(t) = (2, 0, 0)) &= P(\mu_t(x, (2, 0, 0)) = \mu_0) \cdot 1 \\ &\quad + P(\mu_t(x, (2, 0, 0)) = \mu_1) \cdot \frac{\mu_1}{2\mu_0 + \mu_1} \\ &\quad + P(\mu_t(x, (2, 0, 0)) = \mu_2) \cdot \frac{\mu_2}{2\mu_0 + \mu_2} \\ &= \frac{1}{3} \left( 1 + \frac{\mu_1}{2\mu_0 + \mu_1} + \frac{\mu_2}{2\mu_0 + \mu_2} \right). \\ P(\mathbf{M}(t) = (1, 1, 0)) &= P(\mu_t(x, (2, 0, 0)) = \mu_1) \cdot \frac{2\mu_0}{2\mu_0 + \mu_1} \\ &= \frac{1}{3} \cdot \frac{2\mu_0}{2\mu_0 + \mu_1}. \\ P(\mathbf{M}(t) = (1, 0, 1)) &= P(\mu_t(x, (2, 0, 0)) = \mu_2) \cdot \frac{2\mu_0}{2\mu_0 + \mu_2} \\ &= \frac{1}{3} \cdot \frac{2\mu_0}{2\mu_0 + \mu_2}. \end{aligned}$$

Thus

$$\begin{aligned} &P(\mathbf{M}(t) = (2, 0, 0)) + P(\mathbf{M}(t) = (1, 1, 0)) + P(\mathbf{M}(t) = (1, 0, 1)) \\ &= \frac{1}{3} \left( 1 + \frac{\mu_1}{2\mu_0 + \mu_1} + \frac{\mu_2}{2\mu_0 + \mu_2} \right) + \frac{1}{3} \cdot \frac{2\mu_0}{2\mu_0 + \mu_1} + \frac{1}{3} \cdot \frac{2\mu_0}{2\mu_0 + \mu_2} = 1. \end{aligned}$$

The resulting virtual wait at time  $t$  is

$$W(t) = W(t^-) + \mathcal{S}_t = x + \mathcal{S}_t,$$

where  $\mathcal{S}_t$  is the inter start-of-service-depart time, distributed as a mixture

$$\mathcal{S}_t \stackrel{dis}{=} \begin{cases} \text{Exp}_{3\mu_0}, \\ \text{Exp}_{2\mu_0+\mu_1}, \\ \text{Exp}_{2\mu_0+\mu_2}, \end{cases} \text{ with probability } 1/3 \text{ each.}$$

(see Sect. 4.4.1). The sample path will have a jump whose size is distributed as  $\mathcal{S}_t$  at instant  $t$  (see Fig. 4.2).

### 4.5.5 Sample Path of SP Process Revisited

We first describe a typical sample path of the virtual wait in Example 4.4 wherein  $J = 1$  and  $\boldsymbol{\mu} = \{\mu_0, \mu_1\}$ , to facilitate exposition. If  $J > 1$ , sample-path construction would be similar, but with more lines and pages (sheets) in the product space  $\mathbf{T} \times \mathbf{S}$  (see Fig. 4.2). Next, Example 4.5 discusses the general nature of a typical sample path with reference to Example 4.4 and then we outline the mechanics of a *specific* sample path in Example 4.4, based on the M/M/3 queue in Example 4.3 above.

**Example 4.4** Consider M/M/ $c$  with  $c = 3$ ,  $J = 1$ . (Here we take  $J = 1$  for exposition.) A typical sample path of the virtual wait is given in Fig. 4.2).

Arrivals are assigned an exponential service rate from  $\boldsymbol{\mu} = \{\mu_0, \mu_1\}$  with equal probability  $1/2$ . (In general the probabilities can be, e.g.,  $p_0, p_1 = 1 - p_0$ .) The total number of possible configurations is  $\binom{J+c}{J+1} = \binom{4}{2} = 6$ . The full set of configurations is

$$\mathbf{M} = \{00, 10, 01, 11, 20, 02\}.$$

We write  $(2, 0)$  as 20 when  $m_0 = 2, m_1 = 0$ , indicating that 2 servers are occupied with rate  $\mu_0$ ; and similar notation for the other system configurations.

**Example 4.5 General Nature of Sample Path with reference to Example 4.4.** The state space consists of: (1) six discrete points for the zero-wait states  $(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_0 \cup \mathbf{M}_b$ . Thus  $\mathbf{M}_0 = \{00, 10, 01\}$  and  $\mathbf{M}_b = \mathbf{M}_1 = \{11, 20, 02\}$ ; (2) three intervals  $((0, \infty), \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1$ . The three border states are  $(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_b$ .

**Arrival Waits Zero.** Assume an arrival “sees” state  $(0, m)$ ,  $m \in \mathbf{M}_0$ . The SP moves horizontally at time-rate 1 on a line  $\mathbf{T} \times (0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_0$ . If the next arrival occurs before a departure, the SP jumps to a line  $\mathbf{T} \times (0, \mathbf{m}')$ ,  $\mathbf{m}' \in \mathbf{M}_0 \cup \mathbf{M}_b$ , where

$$m'_0 + m'_1 = m_0 + m_1 + 1,$$

because there is one more occupied server. If a departure occurs before an arrival, the SP jumps to a line  $\mathbf{T} \times (0, \mathbf{m}'')$ ,  $\mathbf{m}'' \in \mathbf{M}_0$ , where

$$m''_0 + m''_1 = m_0 + m_1 - 1,$$

because there is one less occupied server. If  $\mathbf{m} = (0, 0)$ , the state can change only due to an arrival.

If an arrival finds the system to be in state  $(0, m)$ ,  $m \in \mathbf{M}_b$  the SP jumps to a sheet  $\mathbf{T} \times ((0, \infty), \mathbf{k})$ ,  $\mathbf{k} \in \mathbf{M}_1$ . Configuration  $\mathbf{k}$  is determined by the service rate assigned to the new arrival, and which server finishes first after the new arrival starts service. Denote the service time of an arrival  $C_{a,t}$  by  $s_t$ . Then  $s_t \stackrel{dis}{=} \text{Exp}_{\mu_t}$  where

$$P(\mu_t = \mu_i) = \frac{1}{2}, i = 0, 1;$$

(see Fig. 4.2).

To fix ideas, let the SP be on the border line  $\mathbf{T} \times (0, 20)$  at arrival instant  $t$ . Thus  $\langle W(t^-), M(t^-) \rangle = \langle 0, 20 \rangle$ .  $C_{a,t}$  starts service upon arrival in the one idle server and is assigned either rate  $\mu_0$  or  $\mu_1$  with probability 1/2 each. Let  $\mathcal{S}_t$  denote the time from the start-of-service of  $C_{a,t}$  until the first departure from the system thereafter. (Since all 3 servers are busy  $\mathcal{S}_t$  is independent of any future arrivals that join the waiting line.)

**Case 1:** Let us assume the service time  $s_t$  has been assigned rate  $\mu_0$ . Then  $\mathcal{S}_t \stackrel{dis}{=} \text{Exp}_{3\mu_0}$  because  $\mathcal{S}_t = \min \{3 \text{ i.i.d. Exp}_{\mu_0}\text{s}\}$ . The SP jumps upward an amount  $\mathcal{S}_t$ . The virtual wait at time  $t$  is

$$W(t) = W(t^-) + \mathcal{S}_t = 0 + \mathcal{S}_t = \mathcal{S}_t.$$

At instant  $t + W(t^-) + \mathcal{S}_t (= t + \mathcal{S}_t)$ , one of the three occupied servers completes service. The service rate of each of the resulting two occupied servers at  $t + \mathcal{S}_t$  must be  $\mu_0$ . By the **look-ahead process**, the configuration at  $t$  is  $M(t) := M(t^-) = 20$ . In this scenario the configuration remains the same as when the test customer arrived. Geometrically, at instant  $t$ , the look-ahead process impels the SP to jump from *line 20* to *page 20*, at a height  $\stackrel{dis}{=}$

$\text{Exp}_{3\mu_0}$  (see Fig. 4.2). In Case 1, the SP jumps from line 20 to page 20 at ordinate  $\underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$ , resulting in  $W(t) \underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$  and  $M(t) = 20$ .

**Case 2.** Let us assume that  $s_t$  has been assigned rate  $\mu_1$ . Then  $\mathcal{S}_t \underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$ . At  $t + \mathcal{S}_t$  one of the three servers completes service. The service rates of the other two still-occupied servers at  $t + \mathcal{S}_t$  are either: (a) both  $\mu_0$  with probability  $\frac{\mu_1}{2\mu_0+\mu_1}$  (the rate- $\mu_1$  server finishes first), or (b)  $\mu_0$  and  $\mu_1$  with probability  $\frac{2\mu_0}{2\mu_0+\mu_1}$  (a rate- $\mu_0$  server finishes first).

In Case 2(a) at instant  $t$ , the SP jumps from line 20 to page 20 at an ordinate  $\underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$ . Thus  $W(t) \underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$  and  $M(t) = 20$ . In Case 2(b) at instant  $t$ , the SP jumps from line 20 to page 11 at ordinate  $\underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$ , resulting in  $W(t) \underset{dis}{=} \text{Exp}_{2\mu_0+\mu_1}$  and  $M(t) = 11$ .

**Arrival Waits a Positive Time.** Assume  $C_{a,t}$  arrives when the state is  $(x, 20)$ ,  $x > 0$  (SP is at ordinate  $x$  on page 20). If the service-rate assignment policy assigns  $s_t \underset{dis}{=} \text{Exp}_{\mu_0}$ , the SP jumps upward an amount  $\text{Exp}_{3\mu_0}$ , and moves with slope  $-1$  steadily on page 20. If the service-rate assignment policy assigns  $s_t \underset{dis}{=} \text{Exp}_{\mu_1}$ , the SP can end up on either page 20 or page 11 at  $t$ . The SP jumps upward to  $W(t) \underset{dis}{=} W(t^-) + \text{Exp}_{2\mu_0+\mu_1}$  and moves with slope  $-1$  steadily on page 20, with probability  $\frac{\mu_1}{2\mu_0+\mu_1}$ . The SP jumps upward to ordinate  $W(t) \underset{dis}{=} W(t^-) + \text{Exp}_{2\mu_0+\mu_1}$  on page 11, with probability  $\frac{2\mu_0}{2\mu_0+\mu_1}$ .

If the SP descends to the bottom of page 20 and hits level 0 from above in a continuous manner before a new arrival occurs, it immediately enters border line 20, and continues its motion along line 20.

### 4.5.6 A Specific Sample Path

We expound further on a possible realization of the SP motion as it traces out the sample path, with reference to Fig. 4.2. Assume that initially the system is empty. The SP moves on line 00. Arrival 1 ( $C_1$ ) sees an empty system. The server-assignment policy assigns  $C_1$  service rate  $\mu_0$ . The SP jumps to, and moves on, line 10.  $C_2$  arrives before  $C_1$  completes service and is also assigned rate  $\mu_0$ . At  $C_2$ 's arrival the SP jumps to line 20.  $C_3$  arrives while both  $C_1$  and

$C_2$  are in service.  $C_3$  receives rate  $\mu_1$ . The SP jumps to an ordinate  $\text{Exp}_{2\mu_0+\mu_1}$ , and if the rate- $\mu_1$  customer finishes first among the three customers in service, the resulting configuration is again 20. The probability of this event is  $\frac{\mu_1}{2\mu_0+\mu_1}$ , due to the memoryless property of exponential variates. This explains why at  $C_3$ 's arrival instant the SP jumps to page 20.

Just before  $C_4$  arrives the SP is descending at slope  $-1$  on page 20.  $C_4$  is assigned service rate  $\mu_0$ . The SP jumps upward an amount  $\text{Exp}_{3\mu_0}$ . It remains on page 20. That is, whichever server finishes first, the two remaining active service rates will be  $\mu_0$ , resulting in configuration 20.  $C_5$  arrives when the SP is on page 20.  $C_5$  is assigned rate  $\mu_1$ . Suppose a server with rate  $\mu_0$  finishes first. The probability of this event is  $\frac{2\mu_0}{2\mu_0+\mu_1}$ . The SP jumps upward by  $\text{Exp}_{2\mu_0+\mu_1}$ . It simultaneously makes a  $20 \rightarrow 11$  transition from page 20 to page 11, since the two remaining occupied servers have rates  $\mu_0$  and  $\mu_1$  when the first service ends. The configuration changes immediately from 20 to 11.

No new arrivals occur prior to the completion of the first rate- $\mu_0$  customer. The SP descends on page 11 with slope  $-1$  and hits level 0 from above, exactly when the first rate- $\mu_0$  customer finishes service. The system now presents a zero wait to a potential arrival. When the SP hits level 0, it enters border line 11 (in Fig. 4.2 it jumps to line 11).  $C_6$  arrives, and starts service immediately.  $C_6$  is assigned rate  $\mu_1$ . The SP jumps to page 02, with probability  $\frac{\mu_0}{\mu_0+2\mu_1}$  ( $\mu_0$ -rate service finishes first), the next configuration is 02.

The system continues to evolve. The SP continues to trace a sample path on the lines and pages according to the probability laws of the model. The sample path gives us a precise picture of the evolving system over time. Construction of the sample path goes hand in hand with understanding the model dynamics, and writing the model equations by inspection.

**Remark 4.7** In Sect. 4.8 below we develop the **steady-state** theory. We will then return to Example 4.3, and formulate the balance equations for the zero-wait probabilities  $P(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M} \equiv \mathbf{M}_0 \cup \mathbf{M}_b$ ; integral equations for the 'partial' pdfs of wait  $f(x, \mathbf{m})$ ,  $x > 0$ ,  $\mathbf{m} \in \mathbf{M}_1$ , and for the total pdf  $\{P_0, f(x), x > 0\}$ .

#### 4.5.7 SP Process Is Markovian

We outline a proof that the SP process is a Markov process. For  $t \geq 0$ , let  $(x, \mathbf{m})_t := \text{event } \{(W(t), M(t)) = (x, \mathbf{m})\}$ . It is required to show that for  $x, y \geq 0$ ,  $\mathbf{m}, \mathbf{k} \in \mathbf{M}$ ,



$$\begin{aligned}
& P((y, \mathbf{k})_{t+h} | (x, \mathbf{m})_t, (W(u), M(u))_{0 \leq u < t}) \\
& = P((y, \mathbf{k})_{t+h} | (x, \mathbf{m})_t), t \geq 0, h > 0.
\end{aligned} \tag{4.9}$$

Formula (4.9) states that the probability of event  $(y, \mathbf{k})_{t+h}$  given that event  $(x, \mathbf{m})_t$  occurred, is independent of the history  $(W(u), M(u))_{0 \leq u < t}$ . We sketch the proof in two steps: (1) **zero-wait states**; (2) **positive-wait states**.

For a Poisson (or non-homogeneous Poisson) process, the probability of more than one event occurring in  $(t, t+h)$  is  $o(h)$  (e.g., Definition 5.3, p. 314, and Definition 5.4, p. 339 in [125]).

### Zero-Wait Non-border States

Assume state  $(0, \mathbf{m})_t \in \left\{ (0, \mathbf{m}) \mid 0 \leq \sum_{j=0}^J m_j \leq c-2 \right\}$  ( $\mathbf{m} \in \mathbf{M}_0, \text{SP} \in S_0$  at time  $t$ ).

**No Departure or Arrival in  $(t, t+h)$**  The state remains  $(0, \mathbf{m})$  in  $(t, t+h)$  iff no arrival or departure occurs during  $(t, t+h)$ , or an event with probability  $o(h)$  occurs. Thus

$$P((x, \mathbf{m})_{t+h} | (x, \mathbf{m})_t) = 1 - \left( \lambda + \sum_{j=0}^J m_j \mu_j \right) h + o(h),$$

which is independent of  $(W(u), M(u))_{0 \leq u < t}$ .

**Arrival in  $(t, t+h)$**  Possibly there is an arrival during  $(t, t+h)$ . The next configuration will have the form

$$\mathbf{m}_{L+} := (m_0, \dots, m_L + 1, \dots, m_J),$$

for some  $L \in \{0, \dots, J\}$ . Then

$$\begin{aligned}
& P((0, \mathbf{m}_{L+})_{t+h} | (0, \mathbf{m})_t) \\
& = (\lambda h + o(h)) \cdot P(\mu_t((0, \mathbf{m})) = \mu_L) \\
& = \lambda h P(\mu_t((0, \mathbf{m})) = \mu_L) + o(h), L \in \{0, \dots, J\}.
\end{aligned} \tag{4.10}$$

Formula (4.10) is the probability that there is an arrival during  $(t, t+h)$  assigned service rate  $\mu_L$ , which is independent of the history given by  $(W(u), M(u))_{0 \leq u < t}$ . Note that

$$\sum_{L=0}^J P(\mu_t((0, \mathbf{m})) = \mu_L) = 1.$$

**Departure in  $(t, t+h)$**  Possibly there is a departure during  $(t, t+h)$ . Let configuration

$$\mathbf{m}_{L-} := (m_0, \dots, \theta_L \cdot (m_L - 1), \dots, m_J), L \in \{0, \dots, J\},$$

where

$$\theta_L = \begin{cases} 1 & \text{if } m_L \geq 1, \\ 0 & \text{if } m_L = 0. \end{cases}$$

Assume  $\mathbf{m} \neq (0, \dots, 0)$ . Then

$$P((0, \mathbf{m}_{L-})_{t+h} | (0, \mathbf{m})_t) = (m_L \cdot \mu_L)h + o(h), \quad (4.11)$$

which is the probability of a rate- $\mu_L$  departure during  $(t, t+h)$  (*rate- $\mu_L$  service finishes first*). Expression (4.11) is independent of the history  $(W(u), M(u)), 0 \leq u < t$ . Note that  $(\sum_{L=0}^J m_L \mu_L)h + o(h)$  is the probability of a departure during  $(t, t+h)$ .

### Zero-Wait Border States

Consider zero-wait border states  $\{(0, \mathbf{m})_t | \sum_{j=0}^J m_j = c - 1\}$  ( $\mathbf{m} \in \mathbf{M}_b, (0, \mathbf{m}) \in \mathbf{S}_b$ ).

**No Arrival in  $(t, t+h)$**  If no arrival or departure occurs, or only a departure occurs, during  $(t, t+h)$ , the Markov property follows similarly as for the zero-wait non-border states given above.

**Arrival in  $(t, t+h)$**  Possibly there is an arrival during  $(t, t+h)$ . In *this case*, the SP jumps to a *positive level on a sheet (page)*. Let configuration

$$\begin{aligned} \mathbf{k} &:= (m_0, \dots, m_L + 1, \dots, m_R - 1, \dots, m_J) \\ &= (k_0, \dots, k_J), \end{aligned}$$

for some  $L, R \in \{0, \dots, J\}$ . Thus  $\sum_{j=0}^J k_j = \sum_{j=0}^J m_j = c - 1$ . Let

$$\nu_L = \sum_{j=0}^J m_j \mu_j + \mu_L.$$

The probability that the SP jumps to sheet  $\mathbf{k}$  during  $(t, t+h)$  and is in state-space interval  $((y, y+dy), \mathbf{k})_{y>0}$  at  $t+h$ , is

$$\begin{aligned} &P((W(t+h), M(t+h)) \in ((y, y+dy), \mathbf{k}) | (0, \mathbf{m})_t) \\ &= (\lambda h + o(h)) \cdot P(\mu_t(0, \mathbf{m}) = \mu_L) \cdot \frac{m_R \mu_R}{\nu_L} \cdot \nu_L \cdot e^{-\nu_L \cdot y} dy \\ &= \lambda h \cdot P(\mu_t(0, \mathbf{m}) = \mu_L) \cdot m_R \mu_R \cdot e^{-\nu_L y} dy + o(h), L \in \{0, \dots, J\}, \end{aligned}$$

which is independent of the history  $(W(u), M(u))_{0 \leq u < t}$ . The right side is the probability that there is an arrival in  $(t, t + h)$ , which is assigned service rate  $\mu_L$ , and a rate- $\mu_R$  service finishes first among the occupied servers, at a time in the state space interval  $(y, y + dy)$ .

### Positive-Wait States

**Arrival In**  $(t, t + h)$  Given  $(x, \mathbf{m})_t$ ,  $x > 0$ , where  $\sum_{j=0}^J m_j = c - 1$ , there may be an arrival during  $(t, t + h)$ . Let

$$\mathbf{k} = (m_0, \dots, m_L + 1, \dots, m_R - 1, \dots, m_J).$$

Reasoning as for zero-wait border states, we obtain

$$\begin{aligned} P((W(t+h), M(t+h)) \in ((x+y, x+y+dy), \mathbf{k}) | (x, \mathbf{m})_t) \\ = \lambda h \cdot P(\mu_t(0, \mathbf{m}) = \mu_L) \cdot m_R \mu_R \cdot e^{-\nu_L(y-x)} dy + o(h), \end{aligned}$$

which is independent of  $(W(u), M(u))_{0 \leq u < t}$ .

### Virtual Wait in $(0, h)$

Consider the case where all servers are occupied, no customers are waiting and  $W(t) \in (0, h)$ , where  $h$  is “small”. Assume a server completes service before a new arrival occurs. Given  $(x, \mathbf{m})_t$ ,  $0 < x < h$ ,  $\sum_{j=0}^J m_j = c - 1$ , we obtain

$$P((0, \mathbf{m})_{t+h} | (x, \mathbf{m})_t) = 1 - \lambda x + o(x).$$

The SP hits level 0 from above in a continuous manner at  $t + x$ . It immediately enters border line  $\mathbf{m}$  corresponding to the border state  $(0, \mathbf{m})$ , and continues its motion in the direction of Time. This is independent of the past history prior to  $t$ .

The above cases cover all possible situations. Formula (4.9) follows in each case, implying that the SP process has the Markov property.

## 4.5.8 Departures from Positive-Wait States

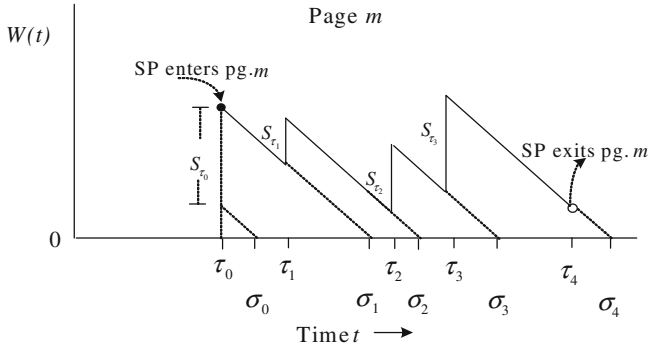
We examine the departure rates during a sojourn on a sheet (page). (Table 4.1 describes the symbols in Fig. 4.3.)

Suppose the SP is at a positive level on page  $\mathbf{m} \in \mathbf{M}_1$  ( $\sum_{j=0}^J m_j = c - 1$  and all  $c$  servers are occupied, including the last arrival). The occupancy number of service rate  $\mu_j$  among the  $c - 1$  servers, not occupied by the last arrival, is  $m_j$ ,  $j \in \{0, \dots, J\}$ .

The single remaining server, which is occupied by the last arrival, may have an arbitrary service rate  $\mu^* \in \mu$ . Assume  $\mu^*$  does not match a positive

**Table 4.1** Description of Symbols in Fig. 4.3

Symbol	Description
$\tau_n$	Arrival instant
$C_{\tau_n}$	Customer that arrives at $\tau_n$
$\sigma_n$	Start of service instant of $C_{\tau_n}$
$\mathcal{S}_{\tau_n}$	$\sigma_{n+1} - \sigma_n =$ inter start-of-service depart time



**Fig. 4.3** SP sojourn on page  $m$ . Departure rate may differ on intervals  $(\tau_0, \sigma_0)$ ,  $(\sigma_0, \sigma_1)$ ,  $(\sigma_1, \sigma_2)$ ,  $(\sigma_2, \sigma_3)$ ,  $(\sigma_3, \tau_4)$ . At instants  $\sigma_0, \sigma_1, \sigma_2$ , arrivals  $C_{\tau_0}, C_{\tau_1}, C_{\tau_2}$  start service. Just after departure instants  $\sigma_0 + \mathcal{S}_{\tau_0}, \sigma_1 + \mathcal{S}_{\tau_1}, \sigma_2 + \mathcal{S}_{\tau_2}$ , the remaining  $c - 1$  servers will have server occupancies  $\mathbf{m} = (m_0, \dots, m_J)$

component in configuration  $\mathbf{m}$ . In order for the SP to *remain* on page  $m$  just after that arrival, the rate- $\mu^*$  server must complete service first among the  $c$  occupied servers (see Fig. 4.3).

While the SP is on page  $m$ , the system exponential departure rate will, in general, differ during inter-departure intervals. These possibly different exponential departure rates have no effect on the *Markov property* of the SP process. The configurations are determined *at arrival instants* (i.e., earlier when service rates are assigned) (Fig. 4.3).

### 4.6 Transient Analysis of Generalized M/M/c

Sections 4.6.1–4.6.6 develop LC relations and definitions leading to the formulation of integro-differential equations for the transient time- $t$  pdf of the virtual wait, in Sect. 4.6.8 below. This development is based on the author’s

working papers [21, 23]. The transient analysis complements the results in [52], which focuses on the generalized M/M/c queue in steady-state. Section 4.6.7 derives the steady-state integral equations by letting  $t \rightarrow \infty$  in the transient equations. Section 4.7 serves as a brief tutorial on writing steady-state model equations using LC and sample paths.

### 4.6.1 Transient PDF of Wait and Downcrossings

We next determine relationships between the transient pdf of wait and sample-path transitions. Let  $\mathcal{D}_t(x, \mathbf{m}) :=$  number of sample-path downcrossings of level  $x$  on page  $\mathbf{m} \in \mathbf{M}_1$  during  $[0, t]$ . Let

$$\mathcal{D}_t(x) = \sum_{\mathbf{m} \in \mathbf{M}_1} \mathcal{D}_t(x, \mathbf{m})$$

denote the total number of downcrossings of level  $x$  on all pages during  $[0, t]$ . Theorem 4.3 connects the instantaneous rate of change of the expected number of downcrossings of level  $x$  in  $[0, t]$ , to the time- $t$  transient pdf of wait at level  $x$ .

**Theorem 4.3** For each configuration  $\mathbf{m} \in \mathbf{M}_1$ ,

$$\frac{\partial}{\partial t} E(\mathcal{D}_t(x, \mathbf{m})) = f_t(x, \mathbf{m}), x > 0, t > 0, \tag{4.12}$$

$$\frac{\partial}{\partial t} E(\mathcal{D}_t(0, \mathbf{m})) = f_t(0^+, \mathbf{m}) (= f_t(0, \mathbf{m})), t > 0, \tag{4.13}$$

$$\frac{\partial}{\partial t} E(\mathcal{D}_t(x)) = f_t(x), x > 0, t > 0, \tag{4.14}$$

$$\frac{\partial}{\partial t} E(\mathcal{D}_t(0)) = f_t(0^+) (= f_t(0)), t > 0. \tag{4.15}$$

**Proof** Fix state-space level  $x > 0$ . Consider instants  $t$  and  $t + h$ , where  $t > 0$ , and  $h > 0$  is small. To prove (4.12) and (4.13) for page  $\mathbf{m}$ , we develop a table similar to (3.10) in Chap. 3 for the M/G/1 queue, and proceed as in the proof of (3.8) and (3.9). Formulas (4.14) and (4.15) follow from the definitions of  $\mathcal{D}_t(x)$  and the total pdf  $f_t(x), x > 0$ . ■

**Corollary 4.1**

$$E(\mathcal{D}_t(x, \mathbf{m})) = \int_{s=0}^t f_s(x, \mathbf{m}) ds, \quad (4.16)$$

$$E(\mathcal{D}_t(0, \mathbf{m})) = \int_{s=0}^t f_s(0^+, \mathbf{m}) ds, \quad (4.17)$$

$$E(\mathcal{D}_t(x)) = \int_{s=0}^t f_s(x) ds, \quad (4.18)$$

$$E(\mathcal{D}_t(0)) = \int_{s=0}^t f_s(0^+) ds. \quad (4.19)$$

**Proof** Integrating both sides of (4.12), (4.13), (4.14) and (4.15) with respect to  $s$  over the interval  $[0, t]$  and applying the initial conditions

$$E(\mathcal{D}_0(x, \mathbf{m})) = E(\mathcal{D}_0(x)) = 0, x \geq 0,$$

yield (4.16), (4.17), (4.18) and (4.19), respectively. ■

**4.6.2 Steady-State PDF of Wait and Downcrossings**

Corollary 4.2 below connects the SP limiting downcrossing rate as  $t \rightarrow \infty$  and the steady-state pdf of wait, at a state-space level. It is analogous to Corollary 3.2 for M/G/1. It also demonstrates the equality of the limit of the instantaneous rate of change of the expected number of downcrossings in  $[0, t]$ , and the limit of the average downcrossing rate over  $[0, t]$ .

Let  $\mathcal{S}_m = ([0, \infty), \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1$ . The results below apply to each page  $\mathbf{T} \times \mathcal{S}_m$ ,  $\mathbf{m} \in \mathbf{M}_1$  as well as to the “book”  $\mathbf{T} \times (\cup_{\mathbf{m} \in \mathbf{M}_1} \mathcal{S}_m)$ .

**Corollary 4.2** *Assume the following limits exist*

$$\lim_{t \rightarrow \infty} f_t(x, \mathbf{m}) \equiv f(x, \mathbf{m}), x \in \mathcal{S}_m, \mathbf{m} \in \mathbf{M}_1.$$

Then

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{D}_t(x, \mathbf{m})) = \lim_{t \rightarrow \infty} \frac{E(\mathcal{D}_t(x, \mathbf{m}))}{t} = f(x, \mathbf{m}), x > 0, \quad (4.20)$$

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{D}_t(0, \mathbf{m})) = \lim_{t \rightarrow \infty} \frac{E(\mathcal{D}_t(0, \mathbf{m}))}{t} = f(0^+, \mathbf{m}) \equiv f(0, \mathbf{m}), \quad (4.21)$$

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{D}_t(x)) = \lim_{t \rightarrow \infty} \frac{E(\mathcal{D}_t(x))}{t} = f(x), x > 0, \tag{4.22}$$

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{D}_t(0)) = \lim_{t \rightarrow \infty} \frac{E(\mathcal{D}_t(0))}{t} = f(0^+) \equiv f(0). \tag{4.23}$$

**Proof** In (4.20), (4.21), (4.22) and (4.23), the equalities of the left-most terms to the pdfs on the right, follow by letting  $t \rightarrow \infty$  in (4.12), (4.13), (4.14) and (4.15), respectively. The equalities of the middle terms to the pdfs on the right, follow by dividing both sides of (4.16), (4.17), (4.18) and (4.19) by  $t > 0$  and letting  $t \rightarrow \infty$ . ■

### 4.6.3 SP $m \rightarrow k$ Transitions

Before discussing the relationship between the transient pdf of wait and SP upcrossings, we define SP  $m \rightarrow k$  transitions. We say that the SP makes an  $m \rightarrow k$  transition at instant  $t_0$  if it *exits* state-space set  $S_m$  and *enters* state-space set  $S_k$  at  $t_0$ . That is, the SP exits  $([0, \infty), m)$  and enters  $([0, \infty), k)$  at  $t_0$ . If  $m = k$ , then an  $m \rightarrow k$  transition maintains the SP on page  $m$  at  $t_0$ . Similar remarks apply to zero-wait lines  $m, k \in M_0$ , or line  $m \in M_b$  and  $S_k$  (see Sects. 2.4.3, 2.4.4 for definitions of entrance and exit).

**$m \rightarrow k$  Upcrossing of a Level** Consider  $S_m, S_k$ . Fix level  $x > 0$ . An  $m \rightarrow k$  upcrossing of level  $x$  occurs at instant  $t_0$  if the SP exits set  $([0, x), m)$  and enters set  $((x, \infty), k)$  at  $t_0$ . That is, the SP makes both an  $m \rightarrow k$  transition and an upcrossing of level  $x$  at  $t_0$ . Thus the SP moves instantaneously (not in Time) from page  $m$  to page  $k$  and from a level below  $x$  to a level above  $x$ . Viewed from the “cover” of the “book”, the upcrossing of level  $x$  resembles an “ordinary” upcrossing of  $x$  by a sample path of the virtual wait in the M/G/1 queue (see Fig. 4.2). Similar definitions apply to line  $m$  and  $S_k$  (page  $k$ ).

**$m \rightarrow k$  Parallel Transition** In some variants of the M/M/c queue, the SP may make “parallel” transitions. The SP makes an  $m \rightarrow k$  parallel transition at  $t_0$  if it exits  $S_m$  from a level  $y$  and enters  $S_k$  at the *same level*  $y$ , at  $t_0$ . SP parallel transitions can also occur in variants of *single-server* queues (e.g., queues with bulk service [20, 93]) and in other stochastic models. The concepts of system configuration, pages (sheets), cover,  $m \rightarrow k$  transitions, etc., are useful in analyzing many other stochastic models.

#### 4.6.4 SP $m \rightarrow k$ Upcrossings Viewed from “Cover”

Let

$$\mathcal{U}_t(x, m, k), m, k \in M_1$$

denote the number of SP  $m \rightarrow k$  upcrossings of level  $x$  during  $[0, t]$ . Denote the *total* number of upcrossings of level  $x$  during  $[0, t]$  (as viewed from the “cover” of the “book”) by

$$\mathcal{U}_t(x) = \sum_{m, k \in M_1} \mathcal{U}_t(x, m, k). \quad (4.24)$$

In (4.24)  $\mathcal{U}_t(x, m, k)$  will be positive only if  $m, k$  are such that page  $k$  is accessible from page  $m$  in one step at an arrival instant (considering lines  $m$  and  $k$  as zero-levels of pages  $m, k$  respectively). For an  $m \rightarrow k$  upcrossing of level  $x$  to occur, the “target” page  $k$  can be either page  $m$  itself ( $k = m$ ) or a different page ( $k \neq m$ ).

#### 4.6.5 Number of Types of $m \rightarrow k$ Upcrossings

A *type* of  $m \rightarrow k$  upcrossing is an ordered pair  $(m, k)$ . The *total* number of possible types of  $m \rightarrow k$  upcrossings depends on how many pages communicate in one step at arrival instants. An upper bound on the total number of possible  $m \rightarrow k$  upcrossings is

$$\begin{aligned} \text{number of ordered pairs } (m, k) &= (\text{number of configurations in } M_1)^2 \\ &= \binom{J+c-1}{c-1}^2 = \binom{J+c-1}{J}^2. \end{aligned}$$

This maximum number  $\binom{J+c-1}{c-1}^2$  is realized only if *all*  $\binom{J+c-1}{c-1}$  pages communicate in one step. In that case, there are  $\binom{J+c-1}{c-1}$  ways to select the “source” page  $m$  and  $\binom{J+c-1}{c-1}$  ways to select the “target” page  $k$  (with replacement).

**Example 4.6** Consider an M/M/c queue with  $c = 3$  and  $J = 1$ , as in Example 4.4 (see Fig. 4.2). The set of configurations corresponding to pages is  $M_1 = \{20, 11, 02\}$ . Here  $\binom{J+c-1}{c-1} = \binom{3}{2} = 3$ . An upper bound on the number of types of  $m \rightarrow k$  transitions (ordered pairs  $(m, k)$ ) is  $3^2 = 9$ . This maximum can be realized only if all configurations in  $M_1$  communicate with each other in one step. This will depend on the probabilities governing the



evolution of the states over time. In the present example, configurations 20 and 02 do not communicate in one step (at an arrival instant). There are *seven* possible **types** of one-step transitions, namely,

$$\{20 \rightarrow 20, 20 \rightarrow 11, 11 \rightarrow 20, 11 \rightarrow 11, 11 \rightarrow 02, 02 \rightarrow 11, 02 \rightarrow 02\}.$$

Transition types  $20 \rightarrow 02$  and  $02 \rightarrow 20$  are not possible.

**The Probability  $p_t(z, \mathbf{m} \rightarrow \mathbf{k})$**  We denote the probability that page  $\mathbf{k}$  is accessible in one step from level  $z$  on page  $\mathbf{m}$  at an arrival instant  $t$ , by  $p_t(z, \mathbf{m} \rightarrow \mathbf{k})$ . Thus for each  $\mathbf{m} \in \mathbf{M}_1$

$$\sum_{\mathbf{k} \in \mathbf{M}_1} p_t(z, \mathbf{m} \rightarrow \mathbf{k}) = 1.$$

Usually, for fixed  $z$ , there is some  $\mathbf{k}$  for which  $p_t(z, \mathbf{m} \rightarrow \mathbf{k}) = 0$ . Then page  $\mathbf{k}$  is not accessible from level  $z$  on page  $\mathbf{m}$  in one step. If such inaccessibility applies for all  $(z, \mathbf{m})$ ,  $z \geq 0$ , then page  $\mathbf{k}$  is not accessible from page  $\mathbf{m}$  in one step. This is the case in Example 4.6: for  $\mathbf{m} = 20$  and  $\mathbf{k} = 02$ ,

$$p_t(z, 20 \rightarrow 02) = p_t(z, 02 \rightarrow 20) = 0, z \geq 0;$$

so, pages  $\mathbf{m}$  and  $\mathbf{k}$  do not communicate in one step.

#### 4.6.6 Transient PDF of Wait and Upcrossings

If a time- $t$  arrival  $C_{a,t}$  finds the state to be  $(z, \mathbf{m})$ , then  $C_{a,t}$  is assigned a service rate  $\mu_t(z, \mathbf{m}) \in \boldsymbol{\mu}$ . We assume that  $\mu_t(z, \mathbf{m})$  is a right continuous with respect to both  $z$  and  $t$ . In Theorem 4.4 we use the fact that  $\mathbf{M}_1 = \mathbf{M}_b = \left\{ \mathbf{m} \mid \sum_{j=0}^J m_j = c - 1 \right\}$  (defined in Sect. 4.4.3).

**Theorem 4.4** For  $\mathbf{m}, \mathbf{k} \in \mathbf{M}_1$ , the instantaneous rate of change of the expected number of  $\mathbf{m} \rightarrow \mathbf{k}$  upcrossings in  $[0, t]$  is given by

$$\begin{aligned} & \frac{\partial}{\partial t} E(\mathcal{U}_t(x, \mathbf{m}, \mathbf{k})) \\ &= \lambda \int_{z=0}^x p_t(z, \mathbf{m} \rightarrow \mathbf{k}) e^{-\nu_t(z, \mathbf{m})(x-z)} dF_t(z, \mathbf{m}), x \geq 0, t \geq 0, \end{aligned} \quad (4.25)$$

where

$$\nu_t(z, \mathbf{m}) = \sum_{j=0}^J m_j \mu_j + \mu_t(z, \mathbf{m}).$$

**Proof** Fix level  $x > 0$  on page  $\mathbf{m}$ , and time  $t > 0$ . Examination of a sample path on page  $\mathbf{m}$  over the time interval  $(t, t + h)$ ,  $h > 0$ , leads to the non-zero values of  $\mathcal{U}_{t+h}(x, \mathbf{m}, \mathbf{k}) - \mathcal{U}_t(x, \mathbf{m}, \mathbf{k})$ , and corresponding probabilities in (4.26) below. We omit  $\mathcal{U}_{t+h}(x, \mathbf{m}, \mathbf{k}) - \mathcal{U}_t(x, \mathbf{m}, \mathbf{k}) = 0$ , which contributes 0 to  $E(\mathcal{U}_{t+h}(x, \mathbf{m}, \mathbf{k}) - \mathcal{U}_t(x, \mathbf{m}, \mathbf{k}))$ . We omit negative values, because  $\{\mathcal{U}_t(x, \mathbf{m}, \mathbf{k})\}_{t \geq 0}$  is a counting process implying  $\mathcal{U}_{t+h}(x, \mathbf{m}, \mathbf{k}) - \mathcal{U}_t(x, \mathbf{m}, \mathbf{k}) \geq 0$ .

$\mathcal{U}_{t+h}(x, \mathbf{m}, \mathbf{k}) - \mathcal{U}_t(x, \mathbf{m}, \mathbf{k})$	Probability
+1	$\lambda h P_0(t) p_t(0, \mathbf{m} \rightarrow \mathbf{k}) e^{-\nu_t(0, \mathbf{m})x} + \lambda h \int_h^x p_t(z, \mathbf{m} \rightarrow \mathbf{k}) e^{-\nu_t(z, \mathbf{m})(x-z)} f_t(z) dz + o(h)$
$\geq 2$	$o(h)$ .

(4.26)

In (4.26), taking the expected value of  $\mathcal{U}_{t+h}(x, \mathbf{m}, \mathbf{k}) - \mathcal{U}_t(x, \mathbf{m}, \mathbf{k})$ , dividing by  $h > 0$  and letting  $h \downarrow 0$ , yields

$$\begin{aligned} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, \mathbf{m}, \mathbf{k})) &= \lambda \cdot P_0(t) \cdot p_t(0, \mathbf{m} \rightarrow \mathbf{k}) \cdot e^{-\nu_t(0, \mathbf{m})x} \\ &+ \lambda \int_{z=0}^x p_t(z, \mathbf{m} \rightarrow \mathbf{k}) \cdot e^{-\nu_t(z, \mathbf{m})(x-z)} \cdot f_t(z) dz \\ &= \lambda \int_{z=0}^x p_t(z, \mathbf{m} \rightarrow \mathbf{k}) \cdot e^{-\nu_t(z, \mathbf{m})(x-z)} \cdot dF_t(z, \mathbf{m}), \end{aligned} \tag{4.27}$$

which is the same as (4.25). ■

**Corollary 4.3** For  $\mathbf{m}, \mathbf{k} \in M_1$ ,

$$\begin{aligned} E(\mathcal{U}_t(x, \mathbf{m}, \mathbf{k})) &= \lambda \int_{s=0}^t \int_{z=0}^x p_s(z, \mathbf{m} \rightarrow \mathbf{k}) \cdot e^{-\nu_s(z, \mathbf{m})(x-z)} \cdot dF_s(z, \mathbf{m}) ds, \quad x \geq 0, t \geq 0. \end{aligned} \tag{4.28}$$

**Proof** In (4.25) change the variable from  $t$  to  $s$  on both sides, integrate with respect to  $s$  over the interval  $[0, t]$ , and apply the initial condition  $E(\mathcal{U}_0(x, \mathbf{m}, \mathbf{k})) = 0$ . This yields (4.28). ■

**Corollary 4.4** Consider the “cover”. For  $x \geq 0, t \geq 0$ ,

$$\frac{\partial}{\partial t} E(\mathcal{U}_t(x)) = \lambda \sum_{m, k \in M_1} \int_{z=0}^x p_t(z, \mathbf{m} \rightarrow \mathbf{k}) \cdot e^{-\nu_t(z, \mathbf{m})(x-z)} \cdot dF_t(z, \mathbf{m}) \quad (4.29)$$

$$\frac{\partial}{\partial t} E(\mathcal{U}_t(0)) = \lambda \sum_{m, k \in M_1} p_t(0, \mathbf{m} \rightarrow \mathbf{k}) \cdot F_t(0, \mathbf{m}). \quad (4.30)$$

**Proof** We define  $\mathcal{U}_t(x), x \geq 0$  in (4.24). Equations (4.29) and (4.30) follow by setting  $x > 0$ , and  $x = 0$ , respectively, in (4.27), and applying (4.24). (The sample path viewed from the cover is the projection of the sample-path segments from all pages onto a single sheet.) ■

**Corollary 4.5** For  $m, k \in M_1$  and  $x \geq 0, t \geq 0$ ,

$$E(\mathcal{U}_t(x)) = \lambda \sum_{m, k} \int_{s=0}^t \int_{z=0}^x p_s(z, \mathbf{m} \rightarrow \mathbf{k}) \cdot e^{-\nu_s(z, \mathbf{m})(x-z)} \cdot dF_s(z, \mathbf{m}) ds,$$

$$E(\mathcal{U}_t(0)) = \lambda \sum_{m, k} \int_{s=0}^t p_s(0, \mathbf{m} \rightarrow \mathbf{k}) \cdot F_s(0, \mathbf{m}) ds.$$

**Proof** In (4.29) and (4.30) change  $t$  to  $s$  and integrate with respect to  $s$  on  $[0, t]$ . Then apply the initial condition  $\mathcal{U}_0(x) = 0, x \geq 0$ . ■

#### 4.6.7 Steady-State PDF of Wait and Upcrossings

Corollary 4.6 below proves

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, \mathbf{m}, \mathbf{k})) = \lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(x, \mathbf{m}, \mathbf{k}))}{t},$$

by relating both limits to the steady-state pdf of wait. Let

$$p(z, \mathbf{m} \rightarrow \mathbf{k}), \nu(z, \mathbf{m}), F(z, \mathbf{m}), \text{ and } f(z, \mathbf{m})$$

be the limiting values of

$$p_t(z, \mathbf{m} \rightarrow \mathbf{k}), \nu_t(z, \mathbf{m}), F_t(z, \mathbf{m}), f_t(z, \mathbf{m}),$$

respectively, as  $t \rightarrow \infty$  (for definition of:  $p_t(z, \mathbf{m} \rightarrow \mathbf{k})$  see Sect. 4.6.5;  $\nu_t(z, \mathbf{m})$  see formula 4.25).

**Corollary 4.6** For  $m, k \in M_1$  and  $x \geq 0$ ,

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, m, k)) &= \lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(x, m, k))}{t} \\ &= \lambda \int_{z=0}^x p(z, m \rightarrow k) \cdot e^{-\nu(z, m)(x-z)} \cdot dF(z, m) \\ &= \lambda p(0, m \rightarrow k) \cdot e^{-\nu(0, m)x} P_0 \\ &\quad + \lambda \int_{z=0}^x p(z, m \rightarrow k) \cdot e^{-\nu(z, m)(x-z)} \cdot f(z, m) dz. \end{aligned} \quad (4.31)$$

**Proof** The equality

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, m, k)) = \lambda \int_{z=0}^x p(z, m \rightarrow k) \cdot e^{-\nu(z, m)(x-z)} \cdot dF(z, m),$$

follows by letting  $t \rightarrow \infty$  on both sides of (4.25). The equality

$$\lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(x, m, k))}{t} = \lambda \int_{z=0}^x p(z, m \rightarrow k) \cdot e^{-\nu(z, m)(x-z)} \cdot dF(z, m)$$

is obtained upon dividing both sides of (4.28) by  $t > 0$ , letting  $t \rightarrow \infty$ , and using L'Hôpital's rule (e.g., Theorem 9, p. 179 in [137]; and many Calculus texts). Equation (4.31) then follows. ■

The next corollary relates the limits

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{U}_t(x)) \text{ and } \lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(x))}{t},$$

for the expected *total number* of upcrossings in  $[0, t]$ , to the steady-state *total* probability distribution of wait.

**Corollary 4.7** For  $x \geq 0$ ,

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\partial}{\partial t} E(\mathcal{U}_t(x)) &= \lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(x))}{t} \\ &= \lambda \sum_{m, k \in M_1} \int_{z=0}^x p(z, m \rightarrow k) \cdot e^{-\nu(z, m)(x-z)} \cdot dF(z, m) \\ &= \lambda \sum_{m, k \in M_1} p(0, m \rightarrow k) \cdot e^{-\nu(0, m)x} P_{0, m} \\ &\quad + \lambda \sum_{m, k \in M_1} \int_{z=0}^x p(z, m \rightarrow k) \cdot e^{-\nu(z, m)(x-z)} \cdot f(z, m) dz. \end{aligned} \quad (4.32)$$

**Proof** The result (4.32) follows from (4.31) and the definition of  $\mathcal{U}_t(x)$  in (4.24). ■

### 4.6.8 Equations for Transient PDF of Wait

We derive the transient model equations for the generalized M/M/c model. These equations comprise: (1)  $\binom{J+c-1}{c-1}$  integro-differential equations for the partial pdfs  $f_t(x, \mathbf{m})$ ,  $x > 0$ ,  $\mathbf{m} \in \mathbf{M}_1$ ; (2)  $\binom{J+c-1}{c-1}$  differential equations for the zero-wait probabilities  $P_t(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1 (= \mathbf{M}_b)$ ; (3)  $\binom{J+c-1}{c-2}$  differential equations for the zero-wait probabilities  $P_t(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_0$ ; (4) one equation for the normalizing condition. Also  $\mathbf{M}_0 = \left\{ \mathbf{m} \mid 0 \leq \sum_{i=0}^J m_i \leq c-2 \right\}$  (see definition in formula 4.8).

We also derive the model equations for the *total* transient mixed pdf of wait  $\{P_0(t), f_t(x)\}_{x>0, t \geq 0}$  (cover of book—see definition in formula 4.5.3).

Formula (4.1) and especially (4.2) of Theorem B (Sect. 4.2.1) play important roles in these derivations. In Theorem B we take the set  $A$  to be an interval in the state space having one of its boundaries equal to  $x$ .

#### Equations for Partial Transient PDFs of Wait

Before stating Theorem 4.5, we introduce/review some definitions.

**Definition 4.5** page  $\mathbf{i} := T \times ((0, \infty), \mathbf{i})$  where system configuration  $\mathbf{i} \in \mathbf{M}_1 (= \mathbf{M}_b)$ —technically page  $\mathbf{i}$  excludes line  $\mathbf{i} := T \times (0, (0, \mathbf{i}))$ , which may be separately depicted, or appended to the bottom of page  $\mathbf{i}$ , in geometric figures (see, e.g., Fig. 4.2);  $J_t^{(0,x)}(\mathbf{k}, \mathbf{m}) :=$  number of  $k \rightarrow m$  jumps that *start* in state set  $((0, x), \mathbf{k})$  during  $[0, t]$  and end in  $((0, \infty), \mathbf{m})$ ;  $J_t^0(\mathbf{k}, \mathbf{m}) :=$  number of  $k \rightarrow m$  jumps that *start* in state set  $(0, \mathbf{k})$  during  $[0, t]$  and end in  $((0, \infty), \mathbf{m})$ ;  $\mathcal{U}_t(x, \mathbf{k}, \mathbf{m}) :=$  number of SP  $\mathbf{k} \rightarrow \mathbf{m}$  transitions that *start* in  $([0, x), k)$  and *jump-upcross* level  $x$  during  $[0, t]$  (start in  $((0, x), \mathbf{k})$  or in  $(0, \mathbf{k})$  during  $[0, t]$  and end in  $((x, \infty), \mathbf{m})$ ).

**Theorem 4.5** (1) The integro-differential equations for  $f_t(x, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1$ , are

$$\begin{aligned} f_t(x, \mathbf{m}) + \lambda \sum_{\mathbf{k} \neq \mathbf{m}} \int_{z=0}^x p_t(z, \mathbf{k}, \mathbf{m}) (1 - e^{-\nu_t(z, \mathbf{m})(x-z)}) f_t(z, \mathbf{k}) dz & \quad (4.33) \\ + \lambda \sum_{\mathbf{k} \neq \mathbf{m}} p_t(0, \mathbf{k}, \mathbf{m}) (1 - e^{-\nu_t(0, \mathbf{k})(x-z)}) P_t(0, \mathbf{k}) & \\ = \frac{\partial}{\partial t} F_t(x, \mathbf{m}) - \frac{\partial}{\partial t} P_t(0, \mathbf{m}) + f_t(0, \mathbf{m}) & \end{aligned}$$

$$\begin{aligned}
 & + \lambda \int_{z=0}^x p_t(z, \mathbf{m}, \mathbf{m}) e^{-\nu_t(z, \mathbf{m})(x-z)} f_t(z, \mathbf{m}) dz \\
 & + \lambda \sum_{\mathbf{k} \neq \mathbf{m}} \int_{z=0}^x p_t(z, \mathbf{m}, \mathbf{k}) f_t(z, \mathbf{m}) dz, \quad x \geq 0, t \geq 0,
 \end{aligned}$$

where configuration  $\mathbf{k} \in \mathbf{M}_1$ .

(2) The differential equation for  $P_t(0, \mathbf{m}), \mathbf{m} \in \mathbf{M}_1$ , is

$$\begin{aligned}
 & f_t(0, \mathbf{m}) + \lambda \sum_{\mathbf{k}} p_t(0, \mathbf{k}, \mathbf{m}) P_t(0, \mathbf{k}) \\
 & = \frac{\partial}{\partial t} P_t(0, \mathbf{m}) + \left( \lambda + \sum_{j=0}^J m_j \mu_j \right) P_t(0, \mathbf{m}) \tag{4.34}
 \end{aligned}$$

where  $\mathbf{k}$  is such that  $\sum_{j=0}^J k_j = c - 2$ .

(3) The differential equations for  $P_t(0, \mathbf{m}), \mathbf{m} \in \mathbf{M}_0$ , are

$$\begin{aligned}
 & \lambda \sum_{\mathbf{r} \neq \mathbf{m}} p_t(0, \mathbf{r}, \mathbf{m}) P_t(0, \mathbf{r}) + \sum_{\mathbf{s} \neq \mathbf{m}} s_j \mu_j p_t(0, \mathbf{s}, \mathbf{m}) P_t(0, \mathbf{s}) \\
 & = \frac{\partial}{\partial t} P_t(0, \mathbf{m}) + \left( \lambda + \sum_{j=0}^J m_j \mu_j \right) P_t(0, \mathbf{m}), \tag{4.35}
 \end{aligned}$$

where state  $(0, \mathbf{m})$  is accessible in one step from state  $(0, \mathbf{r})$  at an arrival instant, and in one step from  $(0, \mathbf{s})$  at a departure instant. That is,

$$\sum_{j=0}^J m_j = \sum_{j=0}^J r_j + 1 = \sum_{j=0}^J s_j - 1.$$

(4) The normalizing condition is

$$\sum_{\mathbf{m} \in \mathbf{M}_0 \cup \mathbf{M}_1} P_t(0, \mathbf{m}) + \sum_{\mathbf{m} \in \mathbf{M}_1} \int_{x=0}^{\infty} f_t(x, \mathbf{m}) dx = 1. \tag{4.36}$$

**Proof (1)** We derive (4.33) by applying Theorem B (Sect. 4.2.1).

**Choose A.** In (4.1) and (4.2), choose  $A := ((0, x), \mathbf{m})$  (i.e.,  $A$  is open interval  $(0, x)$  on page  $\mathbf{m}$ ). The measure of set  $A$  at time  $t$  is

$$P_t(A) = F_t(x, \mathbf{m}) - F_t(0, \mathbf{m}) = F_t(x, \mathbf{m}) - P_t(0, \mathbf{m}).$$

**Entrance rate into A.** The SP can *enter A* by: (i) downcrossing level  $x$  on page  $m$ ; (ii) making a  $k \rightarrow m$  ( $k \neq m$ ) upward jump starting in  $((0, x), k)$ , that **ends in**  $((0, x), m)$ ; (iii) making a jump that *starts from state*  $(0, k)$  ( $k \in M_1$ ) (sometimes located at level 0 on page  $k$  in figures), and **ends in**  $((0, x), m)$ .

The number of SP entrances into set  $A$  during  $[0, t]$  is

$$\begin{aligned} \mathcal{I}_t(A) &= \mathcal{D}_t(x, m) + \sum_{k \neq m \in M_1} J_t^{(0,x)}(k, m) \\ &+ \sum_{k \in M_1} J_t^0(k, m) - \sum_{k \in M_1} \mathcal{U}_t(x, k, m). \end{aligned} \quad (4.37)$$

In (4.37) the algebraic sum

$$\sum_{k \neq m \in M_1} J^{(0,t)}(k, m) + \sum_{k \in M_1} J_t^0(k, m) - \sum_{k \in M_1} \mathcal{U}_t(x, k, m) \quad (4.38)$$

= (number of SP jumps that start in  $([0, x), k)$  on any pages or zero-wait lines  $k \in M_1$ , and end in  $((0, \infty), m)$  – (number of such jumps that end in  $((x, \infty), m)$  on page  $m$  during  $[0, t]$ ). Thus, (4.38) is the number of SP entrances into  $((0, x), m)$  during  $[0, t]$ , due to jumps that start below  $x$  on pages or lines outside of  $T \times ((0, x), m)$  and end in  $((0, x), m)$ . Therefore  $\mathcal{I}_t(A)$  is the *total* number of SP entrances into  $((0, x), m)$  from all sources in one step during  $[0, t]$ .

Taking expected values and then  $\frac{\partial}{\partial t}$  in (4.37) yields

$$\begin{aligned} \frac{\partial}{\partial t} E(\mathcal{I}_t(A)) &= \frac{\partial}{\partial t} E(\mathcal{D}_t(x, m)) + \sum_{k \neq m} \frac{\partial}{\partial t} E(J_t^{(0,x)}(k, m)) \\ &+ \sum_{k \in M_1} \frac{\partial}{\partial t} E(J_t^0(k, m)) - \sum_{k \in M_1} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, k, m)). \end{aligned} \quad (4.39)$$

**Exit Rate of set A.** The SP can *exit set A* by: (i) hitting level 0 on page  $m$  from above in a continuous fashion, (i.e., exiting  $((0, x), m)$  and simultaneously entering state  $(0, m)$ ); (ii) starting in  $((0, x), m)$  at an arrival instant and making an  $m \rightarrow k$  (including  $m \rightarrow m$ ) *upcrossing* of level  $x$ , ending in  $((x, \infty), k)$ ,  $k \in M_1$ ; (iii) starting in  $((0, x), m)$  at an arrival instant and making an  $m \rightarrow k$  ( $k \neq m$ ) jump-transition that ends *below*  $x$  on any page  $k \neq m$ , i.e., in  $((0, x), k)$ ,  $k \in M_1$ ,  $k \neq m$

The total number of exits from set  $A$  during  $[0, t]$  is

$$\begin{aligned} \mathcal{O}_t(A) &= \mathcal{D}_t(0, m) + \sum_{k \in M_1} \mathcal{U}_t(x, m, k) \\ &+ \sum_{k \neq m \in M_1} J_t^{(0,x)}(m, k) - \sum_{k \neq m \in M_1} \mathcal{U}_t(x, m, k). \end{aligned} \tag{4.40}$$

**Explanation of (4.40).** On the right side,  $\mathcal{D}_t(0, m)$  is the number of exits from  $A$  during  $[0, t]$  by downcrossing level 0 on page  $m$  (entering  $(0, m)$ ). The term  $\sum_{k \in M_1} \mathcal{U}_t(x, m, k)$  is the number of SP jump exits from  $([0, x], m)$  during  $[0, t]$  that *upcross* level  $x$  on any page  $k$  (including  $k = m$ ).  $\sum_{k \neq m \in M_1} J_t^{(0,x)}(m, k)$  is the number jump-exits that start in  $((0, x), m)$  and end in  $((0, \infty), k)$  for any  $k \neq m$ . Term  $-\sum_{k \neq m \in M_1} \mathcal{U}_t(x, m, k)$  cancels the extra number of jump-exits from  $([0, x], m)$  during  $[0, t]$  that *upcross* level  $x$  on any page  $k \neq m$ , i.e., ending in  $((x, \infty), k)$ .

Taking expected values and then  $\frac{\partial}{\partial t}$  in (4.40) results in

$$\begin{aligned} &\frac{\partial}{\partial t} E(\mathcal{O}_t(A)) \\ &= \frac{\partial}{\partial t} E(\mathcal{D}_t(0, m)) + \sum_{k \in M_1} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, m, k)) \\ &+ \sum_{k \neq m} \frac{\partial}{\partial t} E\left(\mathcal{U}_t^{(0,x)}(x, m, k)\right) - \sum_{k \neq m} \frac{\partial}{\partial t} E(\mathcal{U}_t(x, m, k)). \end{aligned} \tag{4.41}$$

**Integro-differential Equation:** We substitute in (4.41) from (4.12), (4.13), (4.25). This yields the integro-differential equation (4.33).

(2) We derive (4.34) by letting set  $A = (0, m)$  in Theorem B, and substituting formulas from Section 4.6.1 relating downcrossings and the transient distribution of wait, as in the proof of (1).

(3) We derive (4.35) in a similar manner as in (2).

(4) The final equation is the normalizing condition

$$\sum_{m \in M_0 \cup M_1} P_t(0, m) + \sum_{m \in M_1} \int_{x=0}^{\infty} f_t(x, m) dx = 1.$$

■

**Remark 4.8** In practice we can derive an equivalent set of model equations by letting set  $A = ((x, \infty), m)$ ,  $x > 0$ , in Theorem B (instead of substituting  $((0, x], m)$ ). This choice of  $A$  may simplify the derivation of the model equations for  $f_t(x, m)$ . We would then consider SP jumps that start below and end above level  $x$ . This would yield terms of the form  $e^{-\nu_t(z, m)(x-z)}$



rather than  $(1 - e^{-\nu_t(z, \mathbf{m})(x-z)})$  in the integrands. In real-world applications, writing the integro-differential equations is much simpler than it may seem at this point. Some practice on a few simple models will quickly establish the method. It is very intuitive.

**Remark 4.9** We can generalize the model upon replacing  $\lambda$  by  $\lambda_t$ , depending on  $t$ . The arrival stream would then be a non-homogeneous Poisson process. This generalization holds because the developments in the foregoing sections involving  $\lambda$  are essentially the same if  $\lambda_t$  is substituted for  $\lambda$ .

### Model Equations for Total Transient PDF

In the following theorem, we utilize the previously defined equivalent notation  $F_t(0, \mathbf{m}) \equiv P_t(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1$ ,  $F_t(0) \equiv P_0(t)$ ,  $f_t(0) \equiv f_t(0^+)$ .

**Theorem 4.6** For the total pdf of wait  $\{P_0(t), f_t(x)\}_{x>0}$ , as viewed from the ‘cover’, the following integro-differential and differential equations hold :

$$\begin{aligned} f_t(x) &= \frac{\partial}{\partial t} F_t(x) + \lambda \sum_{\mathbf{m} \in \mathbf{M}_1} \int_{z=0}^x e^{-\nu_t(z, \mathbf{m})(x-z)} dF_t(z, \mathbf{m}) \\ &= \frac{\partial}{\partial t} F_t(x) + \lambda \sum_{\mathbf{m} \in \mathbf{M}_1} P_t(0, \mathbf{m}) e^{-\nu_t(z, \mathbf{m})x} \\ &\quad + \lambda \sum_{\mathbf{m} \in \mathbf{M}_1} \int_{z=0}^x e^{-\nu_t(z, \mathbf{m})(x-z)} f_t(z, \mathbf{m}) dz, \quad x > 0, t \geq 0, \quad (4.42) \end{aligned}$$

$$f_t(0) = \frac{\partial}{\partial t} P_0(t) + \lambda \sum_{\mathbf{m} \in \mathbf{M}_1} P_t(0, \mathbf{m}), \quad t \geq 0. \quad (4.43)$$

**Proof** In Theorem B (Sect. 4.2.1), consider the set

$$\mathbf{A} = (\cup_{\mathbf{m} \in \mathbf{M}_0 \cup \mathbf{M}_1} (0, \mathbf{m})) \cup (\cup_{\mathbf{m} \in \mathbf{M}_1} ((0, x], \mathbf{m}), x > 0).$$

Set  $\mathbf{A}$  includes all  $\binom{J+c}{c-1}$  zero-wait states  $\{(0, \mathbf{m}) | 0 \leq \sum_{j=0}^J m_j \leq c-1\}$ , as well as all positive-wait states  $\{(y, \mathbf{m}) | \sum_{j=0}^J m_j = c-1, y \in (0, x]\}$ .

Every SP *entrance* into  $\mathbf{A}$  must occur from above at level  $x$ . Therefore all entrances are due to (continuous) SP downcrossings of level  $x$ . Every *exit* out of  $\mathbf{A}$  must be due to a jump starting below level  $x$  on a page and ending at a level above level  $x$  on some page. Therefore all SP exits from set  $\mathbf{A}$  are due to upcrossings of level  $x$ .

Thus

$$\begin{aligned}\mathcal{I}_t(A) &= \mathcal{D}_t(x), \mathcal{O}_t(A) = \mathcal{U}_t(x), \\ E(\mathcal{I}_t(A)) &= E(\mathcal{D}_t(x)), E(\mathcal{O}_t(A)) = E(\mathcal{U}_t(x)), \\ \frac{\partial}{\partial t} E(\mathcal{I}_t(A)) &= \frac{\partial}{\partial t} E(\mathcal{D}_t(x)), \frac{\partial}{\partial t} E(\mathcal{O}_t(A)) = \frac{\partial}{\partial t} E(\mathcal{U}_t(x)).\end{aligned}$$

We then substitute these expressions into formulas (4.14), (4.15), (4.29) and (4.30). This substitution yields the integro-differential equation (4.42) and the differential equation (4.43). ■

The normalizing condition

$$P_0(t) + \int_{x=0}^{\infty} f_t(x) dx = 1,$$

is used along with (4.42), (4.43) to solve for the unknown time- $t$  zero-wait probabilities and positive-wait pdfs.

When it is not feasible to obtain an analytical solution, we can use numerical, simulation or approximation techniques to solve for the transient zero-wait probabilities and positive-wait pdfs.

### 4.6.9 Equations for Steady-State PDF of Wait

We obtain the model equations for the steady-state pdf of wait by letting  $t \rightarrow \infty$  in (4.34)–(4.36). All quantities subscripted by  $t$  have limits as  $t \rightarrow \infty$ . We denote the limits utilizing the same notation, omitting subscript  $t$ . If stability holds, then

$$\lim_{t \rightarrow \infty} \frac{\partial}{\partial t} F_t(x, \mathbf{m}) = \lim_{t \rightarrow \infty} \frac{\partial}{\partial t} F_t(0, \mathbf{m}) = 0.$$

This corresponds to the cdf  $F(x, \mathbf{m})$  being independent of  $t$ .

**Theorem 4.7** The integral equation for the steady-state pdf  $f(x, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1$ , is

$$\begin{aligned}f(x, \mathbf{m}) &+ \lambda \sum_{\mathbf{k} \neq \mathbf{m} \in \mathbf{M}_1} \int_{z=0}^x p(z, \mathbf{k}, \mathbf{m}) (1 - e^{-\nu(z, \mathbf{m})(x-z)}) f(z, \mathbf{k}) dz \\ &+ \lambda \sum_{\mathbf{k} \in \mathbf{M}_1} p(0, \mathbf{k}, \mathbf{m}) (1 - e^{-\nu(0, \mathbf{k})(x-z)}) P(0, \mathbf{k}) \\ &= f(0, \mathbf{m}) \\ &+ \lambda \int_{z=0}^x p(z, \mathbf{m}, \mathbf{m}) e^{-\nu(z, \mathbf{m})(x-z)} f(z, \mathbf{m}) dz \\ &+ \lambda \sum_{\mathbf{k} \neq \mathbf{m} \in \mathbf{M}_1} \int_{z=0}^x p(z, \mathbf{m}, \mathbf{k}) f(z, \mathbf{m}) dz, \quad x \geq 0.\end{aligned}\tag{4.44}$$

**Proof** We obtain (4.44) by letting  $t \rightarrow \infty$  in (4.33). ■

**Theorem 4.8** The model equation for the total steady-state pdf is

$$f(x) = \lambda \sum_{\mathbf{m} \in \mathbf{M}_1} P(0, \mathbf{m}) e^{-\nu(z, \mathbf{m})x} + \lambda \sum_{\mathbf{m} \in \mathbf{M}_1} \int_{z=0}^x e^{-\nu(z, \mathbf{m})(x-z)} f(z, \mathbf{m}) dz, \quad x > 0. \quad (4.45)$$

**Proof** Let  $t \rightarrow \infty$  in (4.42). ■

**Remark 4.10** In practice, it is often more efficient to derive balance equations for SP exit/entrance rates with respect to the state-space sets  $((x, \infty), \mathbf{m})_{\mathbf{m} \in \mathbf{M}_1}$ ,  $x > 0$ , rather than with respect to the state-space sets  $((0, x), \mathbf{m})_{\mathbf{m} \in \mathbf{M}_1}$ ,  $x > 0$ . The derived equations will be equivalent, no matter which state-space sets are employed for rate balance.

### Interpretation of Equations in Theorem 4.7 for Sheets

We now interpret (4.44) in terms of rate balance across levels and between pages. This interpretation gives LC power for deriving steady-state model equations by inspecting a typical sample path, in a vast array of complex stochastic models.

In (4.44) the *left* side is the SP *entrance rate* into  $((0, x), \mathbf{m})$ . The term  $f(x, \mathbf{m})$  is the SP downcrossing rate of level  $x$  on page  $\mathbf{m}$ . The term

$$\lambda \sum_{\mathbf{k} \neq \mathbf{m} \in \mathbf{M}_1} \int_{z=0}^x p(z, \mathbf{k}, \mathbf{m}) (1 - e^{-\nu(z, \mathbf{m})(x-z)}) f(z, \mathbf{k}) dz$$

is the rate at which the SP enters composite state  $((0, x), \mathbf{m})$  due to jumps at arrival instants that originate in  $((0, x), \mathbf{k})$  on pages  $\mathbf{k} \neq \mathbf{m}$ . The term

$$\lambda \sum_{\mathbf{k} \in \mathbf{M}_1} p(0, \mathbf{k}, \mathbf{m}) (1 - e^{-\nu(0, \mathbf{k})(x-z)}) P(0, \mathbf{k})$$

is the rate at which the SP enters composite state  $((0, x), \mathbf{m})$  due to jumps that originate at level  $(0, \mathbf{k})$  on any zero-wait line  $\mathbf{k} \in \mathbf{M}_1$ . These three terms exhaust the possible paths by which the SP can enter  $((0, x), \mathbf{m})$ .

The *right* side of (4.44) is the SP *exit rate* of  $((0, x), \mathbf{m})$ . The term  $f(0, \mathbf{m})$  is the rate at which the SP exits  $((0, x), \mathbf{m})$  and simultaneously enters the zero-wait boundary state  $(0, \mathbf{m})$ , due to downcrossings of level 0. The term

$$\lambda \int_{z=0}^x p(z, \mathbf{m}, \mathbf{m}) e^{-\nu(z, \mathbf{m})(x-z)} f(z, \mathbf{m}) dz$$

is the rate at which the SP exits  $((0, x), \mathbf{m})$  and simultaneously enters  $([x, \infty), \mathbf{m})$  due to jumps at arrival instants. The term

$$\lambda \sum_{k \neq \mathbf{m}} \int_{z=0}^x p(z, \mathbf{m}, \mathbf{k}) f(z, \mathbf{m}) dz$$

is the rate at which the SP exits  $((0, x), \mathbf{m})$  and simultaneously enters any page  $\mathbf{k} \neq \mathbf{m}$ . These three terms exhaust the possible paths by which the SP can exit  $((0, x), \mathbf{m})$ .

Thus equation (4.44) is a rate-balance equation of the form:

$$\text{Rate into } ((0, x), \mathbf{m}) = \text{Rate out of } ((0, x), \mathbf{m}),$$

which is a well-known principle for stochastic processes with discrete states, e.g., birth-death processes.

**Interpretation of Equation for Total PDF**

We now provide an LC interpretation of (4.45). We may view the LC analysis of the sheets as a *dissection* of the states of the model (into a partition). The total equation is like a *synthesis*, i.e., reconstruction of the parts into a single whole. This idea helps to derive model equations in complex models directly from sample-path considerations. It utilizes LC ideas for the sheets, lines and the ‘cover’.

In (4.45) the *left* term  $f(x)$  is the total downcrossing rate of level  $x$ , on all pages. On the right side, the term  $\lambda \sum_{\mathbf{m} \in M_1} P(0, \mathbf{m}) e^{-\nu(z, \mathbf{m})x}$  is the total rate at which the SP upcrosses level  $x$  at arrival instants, due to jumps starting at level 0. The term

$$\lambda \sum_{\mathbf{m} \in M_1} \int_{z=0}^x e^{-\nu(z, \mathbf{m})(x-z)} f(z, \mathbf{m}) dz$$

is the total rate at which the SP upcrosses level  $x$  at arrival instants, starting from levels in  $(0, x)$  on all pages  $\mathbf{m} \in M_1$ . We form Eq. (4.45) by rate balance, with respect to level  $x$

$$\text{Downcrossing rate} = \text{Upcrossing rate}.$$

The normalizing condition is

$$P_0 + \int_{x=0}^{\infty} f(x) dx = 1,$$

which too has an LC interpretation. That is, multiply both sides by  $\lambda$ . This yields

$$\lambda P_0 + \lambda \int_{x=0}^{\infty} f(x)dx = \lambda.$$

On the left side,  $\lambda P_0$  is the rate at which the SP makes jumps at arrival instants out of zero-wait states. The term  $\lambda \int_{x=0}^{\infty} f(x)dx$  is the rate at which the SP makes jumps at arrival instants from positive-wait states. The right side  $\lambda$  is the total rate at which the SP makes jumps at arrival instants. The left and right sides are equal.

#### 4.6.10 Discussion of Rate Balance in Complex Models

The rate-balance interpretation provides the analyst with a powerful technique for constructing model equations for steady-state distributions in very complex models. The method is straightforward, intuitive, and relatively easy.

1. Select a state-space interval with boundary  $x$ .
2. Express the SP entrance and exit rates of the interval algebraically in terms of the unknown probability of the interval and/or unknown pdf at  $x$ .
3. Apply rate balance to construct an integral equation (or other type of balance equation) for the probability and/or pdf at  $x$ .
4. Repeat (1)–(3) for every sub-partition of the state space as required, to form a complete system of Volterra integral equations of the second kind (as above), plus other relevant equations, depending on the model.
5. Write the normalizing condition.
6. Solve the entire system of equations simultaneously for the probabilities and pdfs of the model. This can be done analytically, numerically, by approximation, or by LC estimation (see Chap. 9).

**Remark 4.11** The author realized in 1974 that the steady-state model equations discussed here, are really **rate-balance equations**. Originally, these steady-state equations had been derived by starting with Lindley recursions, analogous to those described for M/G/1 in Sect. 1.2 of Chap. 1. The derivation for M/M/c queues started, however, with more complex Lindley recursions.

## 4.7 Example of Steady-State Equations

This Section serves as a brief tutorial on writing steady-state model equations using LC and sample paths. We derive model equations for the steady-state pdf of wait in the *specific* M/M/c queue with  $c = 3$  and  $J = 1$ , discussed in Example 4.4 in Sect. 4.5.5, with a sample path in Fig. 4.2. There are two possible service rates:  $\boldsymbol{\mu} = \{\mu_0, \mu_1\}$ . We make a slight generalization for the *service-rate* assignment policy. For each arrival, the rates  $\{\mu_0, \mu_1\}$  are assigned with probabilities  $\{\alpha_0, \alpha_1\}$ ,  $\alpha_0 + \alpha_1 = 1$  (instead of 1/2 each). Our present example reduces to Example 4.4 if  $\alpha_0 = \alpha_1 = 1/2$ .

We use  $\alpha_0, \alpha_1$  to make it easier to follow the intuitive derivation of the model equations, since  $\alpha_0, \alpha_1$  appear explicitly in the equations.

The set of possible configurations is  $\mathbf{M}_0 \cup \mathbf{M}_1 = \{(m_0, m_1)\}$ , where  $m_j$  denotes the number of servers occupied by customers with service rate  $\mu_j$ ,  $j = 0, 1$ . From the definition of system configuration (Sect. 4.4),

$$0 \leq \sum_{j=0}^1 m_j \leq c - 1 = 2.$$

We abbreviate  $(m_0, m_1)$  as  $m_0m_1$ . There are *six* possible configurations (same as in Example 4.4):

$$\mathbf{M}_0 \cup \mathbf{M}_1 = \{00, 10, 01, 20, 11, 02\}, \quad (4.46)$$

where

$$\mathbf{M}_0 = \{00, 10, 01\}, \quad \mathbf{M}_1 = \{20, 11, 02\}.$$

When an arrival finds more than one server idle, it immediately occupies one of them in accordance with a *server-assignment* rule, and starts service at rate  $\mu_i$  with probability  $\alpha_i$ ,  $i = 0, 1$ .

First we will derive the equations for the zero-wait states (atoms). These are represented in the virtual-wait diagram by the six lines  $\mathbf{T} \times (0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_0 \cup \mathbf{M}_1$  (Fig. 4.2).

Next we will derive the integral equations for the pdfs of the positive-wait states (continuous states). These states are represented by pages  $\mathbf{T} \times ((0, \infty), \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_1$  (Fig. 4.2). Fix level  $x > 0$ . For the equation corresponding to  $\mathbf{m} \in \mathbf{M}_1$ , the left side is the SP exit rate (*out of*) state-space interval  $((x, \infty), \mathbf{m})$ , and the right side is the SP entrance rate (*into*)  $((x, \infty), \mathbf{m})$ . (We use interval  $(x, \infty)$  instead of  $(0, x)$ , since  $(x, \infty)$  results in simpler (equivalent) equations.) Since  $\mathbf{M}_1 = \{20, 11, 02\}$ , there are three pages, three pdfs, and three corresponding integral equations.

**Remark 4.12** To summarize, the zero-wait states are  $(0, m)$ ,  $m \in M_0 \cup M_1$ . The positive-wait states we use for the derivation, are composite states  $((x, \infty), m)$ ,  $m \in M_1$ . We could use alternative state-space intervals having a fixed level- $x$  boundary, such as  $((0, x), m)$  or  $((x, a), m)$ , where constant  $a > x$ , or  $((x, bx), m)$ ,  $b > 1$ , etc. For different interval selections we would derive a different, but equivalent set of model equations. A creative choice of state-space interval may simplify the derivation and final form of the equations. It may lead to new identities or insights about the model. It may also suggest easier ways to obtain solutions of the equations.

The configurations for the zero-wait states are given in  $M_0 \cup M_1$  and for the pages in  $M_1$ , in (4.46) above; (see also Fig. 4.2).

We now derive the model equations; a detailed explanation follows immediately after.

### 4.7.1 Equations for Zero-Wait States

**Notation 4.9** State  $(0, m_0m_1)$  means  $m_0 + m_1$  servers are occupied:  $m_i$  serve at rate  $\mu_i$ ,  $i = 0, 1$ .  $P_{m_0m_1} := P(\text{system is in state } (0, m_0m_1))$ .

Using the principle for discrete states **rate out = rate in**, we obtain the equations for the zero-wait states, as in (4.47). A detailed explanation follows below.

State	Rate out	Rate in
$(0, 00)$	$\lambda P_{00}$	$= \mu_0 P_{10} + \mu_1 P_{01}$
$(0, 10)$	$(\lambda + \mu_0) P_{10}$	$= \lambda \alpha_0 P_{00} + 2\mu_0 P_{20} + \mu_1 P_{11}$
$(0, 01)$	$(\lambda + \mu_1) P_{01}$	$= \lambda \alpha_1 P_{00} + 2\mu_1 P_{02} + \mu_0 P_{11}$
$(0, 20)$	$(\lambda + 2\mu_0) P_{20}$	$= \lambda \alpha_0 P_{10} + f_{20}(0^+)$
$(0, 11)$	$(\lambda + \mu_0 + \mu_1) P_{11}$	$= \lambda \alpha_1 P_{10} + \lambda \alpha_0 P_{01} + f_{11}(0^+)$
$(0, 02)$	$(\lambda + 2\mu_1) P_{02}$	$= \lambda \alpha_1 P_{01} + f_{02}(0^+)$

#### Explanation for Discrete States $(0, m)$ , $m \in M_0 \cup M_b$

In (4.47) the first three equations are derived as in a “bubble” diagram for discrete-state continuous-time Markov chains, using **rate out = rate in**. The last three equations are derived similarly, except for the terms  $f_{20}(0^+)$ ,  $f_{11}(0^+)$ ,  $f_{02}(0^+)$ . These are the exit rates from  $((0, \infty), 20)$ ,  $((0, \infty), 11)$ , and  $((0, \infty), 02)$  into discrete states  $(0, 20)$ ,  $(0, 11)$ ,  $(0, 02)$  respectively. At instants of these exits, the SP simultaneously enters the corresponding line  $T \times (0, 20)$ ,  $T \times (0, 11)$ , or  $T \times (0, 02)$ . It continues its motion.

### 4.7.2 Equations for States $((x, \infty), m)$ , $m \in \mathbf{M}_1$

We now derive the model equations for pages  $m \in \mathbf{M}_1$ . Detailed explanations follow immediately after Eq. (4.51) below.

Rate balance of rates out (left side) and in (right side) for composite state  $((x, \infty), 20)$ ,  $x > 0$ , result in the equation

$$\begin{aligned}
 & f_{20}(x) + \lambda\alpha_1 \frac{2\mu_0}{2\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{20}(y)dy \\
 &= \lambda \left( \alpha_0 e^{-3\mu_0 x} + \alpha_1 \frac{\mu_1}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x} \right) P_{20} \\
 &+ \lambda\alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x} P_{11} + \lambda\alpha_0 \int_{y=0}^x e^{-3\mu_0(x-y)} f_{20}(y)dy \\
 &+ \lambda\alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} \int_{y=0}^x e^{-(2\mu_0 + \mu_1)(x-y)} f_{11}(y)dy \\
 &+ \lambda\alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{11}(y)dy. \tag{4.48}
 \end{aligned}$$

Rate balance for composite state  $((x, \infty), 11)$ ,  $x > 0$ , gives the equation

$$\begin{aligned}
 & f_{11}(x) + \lambda\alpha_1 \frac{\mu_0}{\mu_0 + 2\mu_1} \int_{y=x}^{\infty} f_{11}(y)dy + \lambda\alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{11}(y)dy \\
 &= \lambda \left( \alpha_1 \frac{2\mu_1}{\mu_0 + 2\mu_1} e^{-(\mu_0 + 2\mu_1)x} + \alpha_0 \frac{2\mu_0}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x} \right) P_{11} \\
 &+ \lambda\alpha_1 \frac{2\mu_0}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x} P_{20} + \lambda\alpha_0 \frac{2\mu_1}{\mu_0 + 2\mu_1} e^{-(\mu_0 + 2\mu_1)x} P_{02} \\
 &+ \lambda\alpha_1 \frac{2\mu_1}{\mu_0 + 2\mu_1} \int_{y=0}^x e^{-(\mu_0 + 2\mu_1)(x-y)} f_{11}(y)dy \\
 &+ \lambda\alpha_0 \frac{2\mu_0}{2\mu_0 + \mu_1} \int_{y=0}^x e^{-(2\mu_0 + \mu_1)(x-y)} f_{11}(y)dy \\
 &+ \lambda\alpha_1 \frac{2\mu_0}{2\mu_0 + \mu_1} \int_{y=0}^x e^{-(2\mu_0 + \mu_1)(x-y)} f_{20}(y)dy \\
 &+ \lambda\alpha_0 \frac{2\mu_1}{\mu_0 + 2\mu_1} \int_{y=0}^x e^{-(\mu_0 + 2\mu_1)(x-y)} f_{02}(y)dy \\
 &+ \lambda\alpha_1 \frac{2\mu_0}{2\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{20}(y)dy + \lambda\alpha_0 \frac{2\mu_1}{\mu_0 + 2\mu_1} \int_{y=x}^{\infty} f_{02}(y)dy. \tag{4.49}
 \end{aligned}$$



Rate balance for composite state  $((x, \infty), 02)$ ,  $x > 0$ , gives the equation

$$\begin{aligned}
 & f_{02}(x) + \lambda \alpha_0 \frac{2\mu_1}{\mu_0 + 2\mu_1} \int_{y=x}^{\infty} f_{02}(y) dy \\
 &= \lambda \left( \alpha_1 e^{-3\mu_1 x} + \alpha_0 \frac{\mu_0}{\mu_0 + 2\mu_1} e^{-(\mu_0 + 2\mu_1)x} \right) P_{02} \\
 &+ \lambda \alpha_1 \frac{\mu_0}{\mu_0 + 2\mu_1} e^{-(\mu_0 + 2\mu_1)x} P_{11} + \lambda \alpha_1 \int_{y=0}^x e^{-3\mu_1(x-y)} f_{02}(y) dy \\
 &+ \lambda \alpha_1 \frac{\mu_0}{\mu_0 + 2\mu_1} \int_{y=0}^x e^{-(\mu_0 + 2\mu_1)(x-y)} f_{11}(y) dy \\
 &+ \lambda \alpha_1 \frac{\mu_0}{\mu_0 + 2\mu_1} \int_{y=x}^{\infty} f_{11}(y) dy. \tag{4.50}
 \end{aligned}$$

The normalizing condition is

$$\begin{aligned}
 & P_{00} + P_{10} + P_{01} + P_{20} + P_{11} + P_{02} \\
 &+ \int_{x=0}^{\infty} [f_{20}(x) + f_{11}(x) + f_{02}(x)] dx = 1. \tag{4.51}
 \end{aligned}$$

**Explanation of Equations for States  $((x, \infty), m)$ ,  $m \in M_1$**

**Left Side of Equation (4.48)** In (4.48), the left side represents the SP *exit* rate out of  $((x, \infty), 20)$ . There are two routes by which the SP can exit this composite state: (1) downcrossing level  $x$  on page 20; (2) jumping to page 11 pursuant to an arrival that is assigned rate  $\mu_1$ . The term  $f_{20}(x)$  is the downcrossing rate of level  $x$  on page 20.

The term

$$\lambda \alpha_1 \frac{2\mu_0}{2\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{20}(y) dy$$

is the rate at which the SP jumps to page 11 at arrival instants. In this expression,  $\lambda f_{20}(y) dy$  is the rate at which arrivals find the SP in interval  $(y, y + dy)$  on page 20. The term  $\alpha_1$  is the probability that an arrival gets assigned rate  $\mu_1$ , resulting in two servers having rate  $\mu_0$  and one server having rate  $\mu_1$  just after the arrival starts service. The term  $\frac{2\mu_0}{2\mu_0 + \mu_1}$  is the probability that a rate- $\mu_0$  customer finishes first, causing the SP to jump to page 11. The SP cannot jump to page 02 if an arrival finds the configuration to be 20. The sum of the two terms on the left of side of (4.48) is the *exit* rate of the SP out of  $(x, \infty)$  on page 20.

**Right Side of Equation (4.48)** The right side of (4.48) is the SP entrance rate into  $((x, \infty), 20)$ . The first term

$$\lambda \left( \alpha_0 e^{-3\mu_0 x} + \alpha_1 \frac{\mu_1}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x} \right) P_{20}$$

is the entrance rate into  $((x, \infty), 20)$  due to arrivals that find the state to be  $(0, 20)$ . In it, the product  $\lambda P_{20}$  is the rate at which arrivals find the state to be  $(0, 20)$ . The arrival does not wait, and immediately starts service from the one available server. The term  $\alpha_0 e^{-3\mu_0 x}$  is the product of two probabilities:  $\alpha_0$ , that the arrival is assigned rate  $\mu_0$ ;  $e^{-3\mu_0 x}$ , that the minimum of three independent service times, each having rate  $\mu_0$ , exceeds  $x$ .

The term

$$\alpha_1 \frac{\mu_1}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x}$$

is the product of three probabilities:  $\alpha_1$ , that the arrival is assigned rate  $\mu_1$ ;  $\mu_1 / (2\mu_0 + \mu_1)$ , that the minimum of three service times, two having rate  $\mu_0$  and one having rate  $\mu_1$ , is the rate  $\mu_1$ ;  $e^{-(2\mu_0 + \mu_1)x}$ , that the minimum of the three service times exceeds  $x$ . Both terms result in the SP landing above  $x$  on page 02. The entire term is the rate at which the SP moves from level 0 on page 20 to interval  $(x, \infty)$  on page 20.

The term

$$\lambda \alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} e^{-(2\mu_0 + \mu_1)x} P_{11}$$

is the rate at which arrivals find the state to be  $(0, 11)$  (rate  $\lambda P_{11}$ ), are assigned service rate  $\mu_0$  (probability  $\alpha_0$ ), the minimum service time is a rate- $\mu_1$  service (probability  $\mu_1 / (2\mu_0 + \mu_1)$ ), and the minimum exceeds  $x$  (probability  $e^{-(2\mu_0 + \mu_1)x}$ ). This is the rate at which the SP moves from discrete level 0 on page 11 to  $(x, \infty)$  on page 20.

The term

$$\lambda \alpha_0 \int_{y=0}^x e^{-3\mu_0(x-y)} f_{20}(y) dy$$

is the rate at which arrivals find the state to be  $(y, 20)$ ,  $y \in (0, x)$ , are assigned service rate  $\mu_0$  (probability  $\alpha_0$ ), and the minimum of three service times each having rate  $\mu_0$  exceeds  $x - y$  (probability  $e^{-3\mu_0(x-y)}$ ) integrated over all  $y \in (0, x)$ . This is the rate at which the SP moves from  $(0, x)$  on page 20 to  $(x, \infty)$  on page 20 (makes  $20 \rightarrow 20$  upcrossings of  $x$ ).

The term

$$\lambda \alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} \int_{y=0}^x e^{-(2\mu_0 + \mu_1)(x-y)} f_{11}(y) dy$$

is the rate at which arrivals find the state to be in  $((y, y + dy), 11)$ ,  $y \in (0, x)$  (factor  $\lambda f_{11}(y)dy$ ), are assigned service rate  $\mu_0$ , the rate- $\mu_1$  service ends first, and the minimum of three exponential r.v.s (two having rate  $\mu_0$  and one rate  $\mu_1$ ) exceeds  $x - y$ , integrated over all  $y \in (0, x)$ . This is the rate at which the SP moves from  $(0, x)$  on page 11 to  $(x, \infty)$  on page 20 (makes  $11 \rightarrow 20$  upcrossing of  $x$ ).

The term

$$\lambda \alpha_0 \frac{\mu_1}{2\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{11}(y)dy$$

is the rate at which arrivals find the state to be in  $((y, y + dy), 11)$ ,  $y > x$ , are assigned service rate  $\mu_0$ , the rate- $\mu_1$  service finishes first, and the minimum of three exponential service times (two having rate  $\mu_0$  and one having rate  $\mu_1$ ) has any value in  $(x, \infty)$ . This is the rate at which the SP moves from  $(x, \infty)$  on page 11 to  $(x, \infty)$  on page 20 (makes  $11 \rightarrow 20$  transition, from and to, points above  $x$ ).

**Integral Equations (4.49) and (4.50)**

We derive integral equations (4.49) and (4.50) for the pdfs  $f_{11}(x)$  and  $f_{02}(x)$  (pages 11 and 02), in a similar manner as for  $f_{20}(x)$  above.

**Normalizing Condition**

The normalizing condition (4.51) ensures that the sum of all zero-wait and positive-wait probabilities is 1.

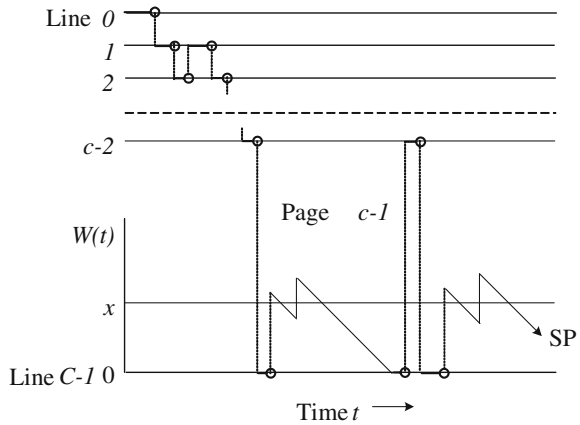
**4.8 Standard M/M/c: Steady-State Analysis**

We analyze the standard M/M/c queue as a special case of the generalized M/M/c queue developed in Sects. 4.3–4.7. It is instructive to derive known results for M/M/c using LC. *The standard M/M/c queue is one of the first models the author analyzed in 1974, to validate the LC method* (see pp. 37–39 in [11]).

We assume the number of servers is  $c \geq 2$ , each customer receives the same exponential service rate  $\mu$ , and  $\lambda < c\mu$ . Using the notation of Sect. 4.3, we have here  $J = 0$ ,  $\boldsymbol{\mu} = \{\mu_0\} := \{\mu\}$ . A system configuration has one component  $m_0$ , which can take values in  $\{0, 1, \dots, c - 1\}$ . The virtual wait process is denoted as  $\{W(t)\}_{t \geq 0}$ .

In this model, a system configuration is a scalar  $m_0 :=$  *number of customers in the other servers just after an arrival starts service*. Thus  $m_0 \in \{0, 1, \dots, c - 1\}$ . Equivalently  $m_0$  is the number of *other* occupied servers at a start of service instant. The set of all configurations,  $\mathbf{M} = \mathbf{M}_0 \cup \mathbf{M}_1$ , has size  $\binom{J+c}{c-1} = \binom{0+c}{c-1} = \binom{c}{1} = c$  (see Sect. 4.4.2). That is,

**Fig. 4.4** Sample path of  $\{W(t)\}_{t \geq 0}$  in standard M/M/c queue. There are  $c$  lines and one page. Line  $c - 1$  is at level 0 of page  $c - 1$



$$M_0 = \{0, 1, \dots, c - 2\}, \quad M_1 = \{c - 1\}.$$

(Recall that  $M_1 = M_b$ , the set of ‘border’ configurations.)

A sample path of  $\{W(t)\}_{t \geq 0}$  has  $c$  lines for the zero-wait states  $(0, j)$ ,  $j = 0, \dots, c - 1$ , and one page (sheet) for the composite state  $((0, \infty), c - 1)$  (Fig. 4.4). Line  $c - 1$ , the border line corresponding to state  $(0, c - 1)$ , is usually placed as the bottom line of page  $c - 1$ , but is arbitrarily located among the other 0-wait states in Fig. 4.4. This does not change the analysis because rate balance across level  $x > 0$  (downcrossing rate = upcrossing rate) is equivalent to rate balance between sets  $((x, \infty), c - 1)$  and  $([0, x], c - 1)$  i.e., rate from  $((x, \infty), c - 1)$  into  $([0, x], c - 1)$  = rate from  $([0, x], c - 1)$  into  $((x, \infty), c - 1)$ .

Denote the zero-wait probabilities as  $P_n$ ,  $n = 0, \dots, c - 1$ , the pdf of wait as  $f(x)$ ,  $x > 0$ , and the steady-state cdf of wait by  $F(x)$ ,  $x \geq 0$ . Then

$$F(x) = \sum_{n=0}^{c-1} P_n + \int_0^x f(x)dx, \quad x \geq 0,$$

$$F(0) = \sum_{n=0}^{c-1} P_n.$$

### 4.8.1 Equations for Steady-State PDF of Wait

We derive model equations for the steady-state pdf of wait, and give further explanations in Sect. 4.8.2 below.

**Zero-Wait States**

For the zero-wait states (atoms) the model equations are (using *rate out = rate in*)

$$\begin{aligned} \lambda P_0 &= \mu P_1 \\ (\lambda + \mu) P_1 &= \lambda P_0 + 2\mu P_2 \\ (\lambda + 2\mu) P_2 &= \lambda P_1 + 3\mu P_3 \\ &\dots \\ (\lambda + (c - 2)\mu) P_{c-2} &= \lambda P_{c-3} + (c - 1)\mu P_{c-1} \\ (\lambda + (c - 1)\mu) P_{c-1} &= \lambda P_{c-2} + f(0^+). \end{aligned} \tag{4.52}$$

The term  $f(0^+)$  in the last equation in (4.52) connects the *pdf of a continuous random variable (waiting time)* with the *probabilities of atoms (states (0, c - 1) and (0, c - 2))*. This observation (and other examples) led the author to coin the term “*border state*” (i.e., state (0, c - 1) in the present context).

**Positive-Wait States**

For the composite state ((0, ∞), c - 1) (the single page) the model equation is

$$f(x) = \lambda P_{c-1} e^{-c\mu x} + \lambda \int_{y=0}^x e^{-c\mu(x-y)} f(y) dy, \quad x > 0. \tag{4.53}$$

Composite state ((0, ∞), c - 1) is accessible in one step at an arrival instant, only from the border state (0, c - 1). The normalizing condition is

$$F(\infty) = \sum_{n=0}^{c-1} P_n + \int_{y=0}^{\infty} f(x) dx = 1. \tag{4.54}$$

**4.8.2 Explanation of Equations (4.52) and (4.53)**

**Linear Equations (4.52)**

Equation (4.52) are rate-balance equations, which equate SP rates out of, and into, the discrete zero-wait states (0, n), n = 0, . . . , c - 1. The term  $f(0^+)$  (:=  $f(0)$ ) is the SP downcrossing rate of level 0, i.e., the SP *entrance rate into state (0, c - 1) from above*. (In sample-paths, line c - 1 may be equally placed at level 0 of page (c - 1). If it is placed separately as in Fig. 4.4, we can still imagine it to be *at level 0* of the page with respect to SP motion.)

**Integral Equation (4.53)**

To derive the positive-wait integral equation (4.53) consider composite state ((x, ∞), c - 1) on the (single) page (Fig. 4.4). We equate the SP *exit rate*

(i.e., downcrossing rate of level  $x$ ) to the *entrance rate* (i.e., ‘upcrossing’ rate of level  $x$  starting from line  $c - 1$  thought of as being at the bottom of the page). The downcrossing rate of level  $x$  is  $f(x)$  (see Corollary 4.2 in Sect. 4.6.2).

The SP entrance rate into  $((x, \infty), c - 1)$  is from two sources:

(1) Entrances are generated by jumps due to arrivals when the state is the border state  $(0, c - 1)$ , starting from level 0 of the page and ending above level  $x$  on the page. Since there is only one page, the only access to  $((x, \infty), c - 1)$  in one step from a zero-wait state is from state  $(0, (c - 1))$ , i.e., line  $c - 1$  in the sample path. The SP entrance rate from this source is  $\lambda P_{c-1} \cdot P(S > x)$ , where  $P_{c-1}$  is the limiting probability of state  $(0, c - 1)$ , and  $S$  is the *inter start-of-service depart time*. (See Sect. 4.4.1 for a discussion of *inter start-of-service depart time*.) Random variable  $S \stackrel{dis}{=} \text{Exp}_{c\mu}$ , since there would be  $c$  customers with rate  $\mu$  in service just after such an arrival starts service, and  $S := \text{minimum of } c \text{ i.i.d. } \text{Exp}_{\mu} \text{ random variables}$ . Thus,  $P(S > x) = e^{-c\mu x}$ . This gives the term  $\lambda P_{c-1} e^{-c\mu x}$  in (4.53).

(2) Entrances into  $((x, \infty), c - 1)$  are generated by jumps due to arrivals when the state is a continuous state  $(y, c - 1)$ ,  $y \in (0, x)$ . Such jumps start at level  $y$  and end above level  $x$ . Just after such an arrival begins service ( $y$  after its arrival), all  $c$  servers will be occupied and  $S \stackrel{dis}{=} \text{Exp}_{c\mu}$ , independent of any new arrivals to the system. The SP will enter  $((x, \infty), c - 1)$  with probability  $e^{-c\mu(x-y)}$ . This leads to Eq. (4.53).

### 4.8.3 Solution of Equations

We first solve (4.53). Differentiating both sides with respect to  $x$  and solving the resulting first-order differential equation, gives

$$f(x) = A e^{-(c\mu-\lambda)x}, x > 0,$$

where  $A$  is a constant. Letting  $x \downarrow 0$ , we get the initial condition

$$f(0) = A = \lambda P_{c-1} \tag{4.55}$$

since  $f(0)$  ( $:= f(0^+)$ ) is the SP downcrossing rate of level 0, and  $\lambda P_{c-1}$  is the ‘upcrossing’ rate of level 0 (rate of egress from  $(0, c - 1)$  above). (Equivalently,  $f(0)$  is the exit rate out of  $((0, \infty), c - 1)$  and  $\lambda P_{c-1}$  is the entrance rate into  $((0, \infty), c - 1)$ .) Thus  $A = \lambda P_{c-1}$  and

$$f(x) = \lambda P_{c-1} e^{-(c\mu-\lambda)x}, x > 0. \tag{4.56}$$

Note that the condition (4.55) is itself a rate-balance equation for the rates out of, and into,  $((0, \infty), c - 1)$ .

Next, from (4.52) and (4.56) we obtain

$$\begin{aligned} P_n &= \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} P_0, n = 0, \dots, c - 1, \\ P_{c-1} &= \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!} P_0. \end{aligned} \quad (4.57)$$

Substituting (4.57) into (4.56) gives

$$f(x) = \lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!} P_0 \cdot e^{-(c\mu-\lambda)x}, x > 0.$$

The normalizing condition (4.54) is

$$\left(\sum_{n=0}^{c-1} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!}\right) P_0 + \lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!} P_0 \int_{x=0}^{\infty} e^{-(c\mu-\lambda)x} dx = 1.$$

This gives the well-known value

$$P_0 = \frac{1}{\sum_{n=0}^{c-1} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} + \left(\frac{\lambda}{\mu}\right)^c \frac{c\mu}{c!(c\mu-\lambda)}}. \quad (4.58)$$

The cdf of wait is

$$\begin{aligned} F(x) &= P_0 + \int_{y=0}^x \lambda P_{c-1} e^{-(c\mu-\lambda)y} dy \\ &= P_0 \left( 1 + \lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!(c\mu-\lambda)} \left( 1 - e^{-(c\mu-\lambda)x} \right) \right), x \geq 0. \end{aligned} \quad (4.59)$$

### Boundedness of PDF of Wait

From (4.56)  $f(x) < \lambda, x > 0$ , since  $P_{c-1} < 1$  and  $e^{-(c\mu-\lambda)x} < 1$  ( $c\mu - \lambda > 0$ ).

#### 4.8.4 CDF and PDF of Wait Geometrically

It is insightful and intuitive to derive the steady-state cdf and pdf of wait geometrically, directly from sample path properties. This derivation bypasses model equation (4.53). A similar geometric derivation for the cdf of wait in the M/M/1 queue is given in Sect. 3.5.6.

Consider level  $x > 0$  on the single page (Fig. 4.4). Rate balance across level  $x$  applies the principle

$$\text{Upcrossing rate of } x = \text{Downcrossing rate of } x = f(x).$$

Equivalently, in symbols

$$\lim_{t \rightarrow \infty} \frac{\mathcal{U}_t(x)}{t} = \lim_{t \rightarrow \infty} \frac{\mathcal{D}_t(x)}{t} = f(x) \text{ (a.s.)},$$

or

$$\lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(x))}{t} = \lim_{t \rightarrow \infty} \frac{E(\mathcal{D}_t(x))}{t} = f(x).$$

The sojourn time above level  $x > 0$  on the page, initiated by each upcrossing of  $x$ , is := *busy period of a standard*  $M_\lambda/M_{c\mu}/1$  queue with arrival rate  $\lambda$  and *service rate*  $c\mu$  because, when the SP is on the page, all  $c$  servers are occupied and each is serving at rate  $\mu$ . Thus the inter start-of-service depart time (see Definition 4.3 in Sect. 4.4.1)  $\mathcal{S}_{dis}$  = *size of each jump ending on the page* =  $\text{Exp}_{c\mu}$ . Moreover, by the memoryless property, excess jumps above level  $x$  are =  $\text{Exp}_{c\mu}$ .

Let  $a_x$  denote an SP sojourn time above  $x$ . Then  $a_x$  = *busy period in*  $M_\lambda/M_{c\mu}/1$  ( $\lambda < c\mu$ ). Thus

$$E(a_x) = \frac{1}{c\mu - \lambda}, \quad (4.60)$$

*independent* of  $x$ , since the expected value of the busy period in  $M_\lambda/M_{c\mu}/1$  is  $1/(c\mu - \lambda)$ .

Let  $d_x$  := *inter-downcrossing time at level*  $x \geq 0$ . Since level- $x$  downcrossings are regenerative points, similarly as in Sect. 3.4.15 we have

$$E(d_x) = 1/f(x). \quad (4.61)$$

The renewal reward theorem (Sect. 3.79), now yields



$$\frac{E(a_x)}{E(d_x)} = \lim_{t \rightarrow \infty} \frac{\text{time } W(\cdot) \in (x, \infty) \text{ during } (0, t)}{t} = 1 - F(x),$$

$$\frac{1/(c\mu - \lambda)}{1/f(x)} = 1 - F(x),$$

or

$$\frac{f(x)}{1 - F(x)} = c\mu - \lambda, \quad x > 0, \quad (4.62)$$

equivalent to the differential equation

$$\frac{\frac{d}{dx}(1 - F(x))}{1 - F(x)} = -(c\mu - \lambda),$$

$$\frac{d}{dx} \ln(1 - F(x)) = -(c\mu - \lambda),$$

with solution

$$1 - F(x) = A \cdot e^{-(c\mu - \lambda)x},$$

where  $A$  is a constant, evaluated by letting  $x \downarrow 0$ , and yielding the cdf of wait

$$F(x) = 1 - (1 - F(0))e^{-(c\mu - \lambda)x}, \quad x \geq 0, \quad (4.63)$$

where  $F(0) = P(\text{zero wait})$ . Taking  $dF(x)/dx$ ,  $x > 0$ , in (4.63) gives the pdf of wait

$$f(x) = (1 - F(0))(c\mu - \lambda)e^{-(c\mu - \lambda)x}, \quad x > 0. \quad (4.64)$$

We next employ the equations in (4.57) to get

$$F(0) = \sum_{n=0}^{c-1} P_n = P_0 \sum_{n=0}^{c-1} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!}. \quad (4.65)$$

Note that  $f(0) = \lambda P_{c-1}$ , i.e., the SP entrance rate into state  $(0, c - 1)$  from above (downcrossing rate of level 0) is equal to the SP exit rate from state  $(0, c - 1)$  at arrival instants. Letting  $x \downarrow 0$  in (4.64) yields

$$f(0) = (1 - F(0))(c\mu - \lambda) = \lambda P_{c-1}. \quad (4.66)$$

From (4.66) and (4.57)

$$F(0) = 1 - \frac{\lambda}{c\mu - \lambda} P_{c-1} = 1 - \frac{\lambda}{c\mu - \lambda} \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!} P_0. \quad (4.67)$$

Substituting the value of  $F(0)$  from (4.65) into (4.67) and solving for  $P_0$  gives (4.58). The upshot is two different ways to determine  $P_0$ ; and two different, equivalent formulas for  $f(x)$ ,  $x > 0$ : (4.56) and (4.64).

**Remark 4.13** Another way to obtain the second equality in (4.66) is to note that the SP expected sojourn time above 0 is

$$E(a_0) = E(\text{busy period of } M_\lambda/M_{c\mu}/1) = \frac{1}{c\mu - \lambda}.$$

The proportion of time the SP spends above level 0 is therefore

$$\lim_{t \rightarrow \infty} \frac{E(\mathcal{U}_t(0))}{t} \cdot \frac{1}{c\mu - \lambda} = \lambda P_{c-1} \cdot \frac{1}{c\mu - \lambda} = 1 - F(0).$$

**Busy Period in M/M/c**

Note that  $a_0$  is equal to a busy period in M/M/c, denoted by  $\mathcal{B}_{c-1,c}$ , defined as the time measured from an arrival instant when the state is  $(0, c - 1)$  until the first departure instant thereafter that leaves the system in state  $(0, c - 1)$  again. (The arrival increases the number in the system to  $c$ . The departure decreases the number to  $c - 1$ .) Since  $a_x \equiv a_0$ ,  $x \geq 0$ ,

$$E(\mathcal{B}_{c-1,c}) = E(a_0) = E(a_x) = \frac{1}{c\mu - \lambda}, x \geq 0. \tag{4.68}$$

We also call  $\mathcal{B}_{c-1,c}$  a  $[c - 1, c]$  busy period.

**4.8.5 PMF of Number in the System**

We use the foregoing pdf of wait (4.56) to derive  $P_n$ ,  $n = c, c + 1, \dots$ . This approach is the reverse order of the usual derivation, which first derives the pmf (probability mass function) of the number-in-system using a birth-death analysis. It then obtains the pdf of wait by conditioning on the number in the system when there is an arrival. The method we apply here utilizes partly birth-death analysis and partly LC. It provides a different perspective on the M/M/c model.

Due to Poisson arrivals,  $P_n = a_n = d_n$ , where  $a_n, d_n$  are the steady-state probabilities of  $n$  units in the system just before an arrival, and just after a departure, respectively (in this Section). Reasoning as for M/M/1 (see Sect. 3.5.3), we get

$$P_n = d_n = P(n - c \text{ arrivals during a waiting time}), n = c, c + 1, \dots$$

Substituting from (4.56) and (4.57)

$$\begin{aligned}
 P_n &= \int_{x=0}^{\infty} \frac{e^{-\lambda x} (\lambda x)^{n-c}}{(n-c)!} f(x) dx \\
 &= \left(\frac{\lambda}{\mu}\right)^{n-c+1} \frac{1}{c^{n-c+1}} P_{c-1} \int_{x=0}^{\infty} c\mu e^{-c\mu x} \frac{(c\mu x)^{n-c}}{(n-c)!} dx \\
 &= \left(\frac{\lambda}{\mu}\right)^n \frac{1}{c^{n-c} c!} P_0, n = c, c + 1, \dots .
 \end{aligned}$$

In summary, we obtain the well-known formulas (e.g., p. 67 in [84])

$$\left. \begin{aligned}
 P_0 &= \frac{1}{\sum_{n=0}^{c-1} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} + \left(\frac{\lambda}{\mu}\right)^c \frac{c\mu}{c!(c\mu-\lambda)}} \\
 P_n &= \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} P_0, n = 0, \dots, c - 1, \\
 P_n &= \left(\frac{\lambda}{\mu}\right)^n \frac{1}{c^{n-c} c!} P_0, n = c, c + 1, \dots .
 \end{aligned} \right\} \quad (4.69)$$

The probability that all servers are occupied is

$$\begin{aligned}
 \sum_{n=c}^{\infty} P_n &= P(\text{wait} > 0) = \int_{x=0}^{\infty} f(x) dx \\
 &= \lambda P_{c-1} \int_{x=0}^{\infty} e^{-(c\mu-\lambda)x} dx = \frac{\lambda}{c\mu - \lambda} P_{c-1} \\
 &= \frac{\lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{c!}}{c\mu - \lambda} P_0.
 \end{aligned} \quad (4.70)$$

The probability that there is at least one idle server is

$$\sum_{n=0}^{c-1} P_n = P(\text{wait} = 0) = 1 - \frac{\lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{c!}}{c\mu - \lambda} P_0. \quad (4.71)$$

### 4.8.6 Inter-downcrossing and Sojourn Times

Consider  $d_x, a_x, b_x$  ( $x \geq 0$ ), respectively: time between successive SP downcrossings of level  $x$ ; sojourn time above  $x$  initiated by an upcrossing of  $x$ ; sojourn time at or below  $x$  initiated by a downcrossing of  $x$ . Formula (4.61)

shows

$$E(d_x) = \frac{1}{f(x)} = \frac{e^{(c\mu-\lambda)x}}{\lambda P_{c-1}}, x \geq 0; \quad (4.72)$$

formula (4.60) shows that, independent of  $x$ ,

$$E(a_x) = \frac{1}{c\mu - \lambda}.$$

Note that

$$\lim_{t \rightarrow \infty} \frac{(\text{time that the SP is above } x \text{ during } (0, t))}{t} = 1 - F(x);$$

by the renewal reward theorem (Sect. 3.4.9),

$$\begin{aligned} \frac{E(a_x)}{E(d_x)} &= \frac{E(a_x)}{1/f(x)} = 1 - F(x), \\ E(a_x) &= \frac{1 - F(x)}{f(x)} = \frac{1}{c\mu - \lambda}, x > 0, \end{aligned}$$

The last equality above corroborates formula (4.62) when solving for  $f(x)$  geometrically in Sect. 4.8.4. Also, we can validate (4.62) in Sect. 4.8.4 using

$$\begin{aligned} \frac{1 - F(x)}{f(x)} &= \frac{\int_{y=x}^{\infty} f(y) dy}{\lambda P_{c-1} e^{-(c\mu-\lambda)x}} \\ &= \frac{\int_{y=x}^{\infty} \lambda P_{c-1} e^{-(c\mu-\lambda)y} dy}{\lambda P_{c-1} e^{-(c\mu-\lambda)x}} = \frac{1}{c\mu - \lambda}. \end{aligned}$$

Note that  $F(x) = \lim_{t \rightarrow \infty} (\text{time the SP spends at or below } x \text{ during } (0, t))/t$ . Each instant that the SP downcrosses  $x \geq 0$  is a regenerative point, due to the memoryless property of the interarrival times. From the renewal reward theorem (i.e., the theory of regenerative processes, e.g., [134])  $E(b_x)/E(d_x) = F(x)$ , implying

$$\begin{aligned} E(b_x) &= \frac{F(x)}{f(x)} = \frac{1 - (1 - F(0))e^{-(c\mu-\lambda)x}}{\lambda P_{c-1} e^{-(c\mu-\lambda)x}} \\ &= \frac{e^{(c\mu-\lambda)x}}{\lambda P_{c-1}} - \frac{(1 - F(0))}{\lambda P_{c-1}} \\ &= \frac{1}{\lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!} P_0} \left( e^{(c\mu-\lambda)x} - (1 - P_0) \sum_{n=0}^{c-1} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!} \right) \end{aligned}$$

$$= \frac{e^{(c\mu-\lambda)x} - 1}{\lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!} P_0} + \frac{\sum_{n=0}^{c-1} \left(\frac{\lambda}{\mu}\right)^n \frac{1}{n!}}{\lambda \left(\frac{\lambda}{\mu}\right)^{c-1} \frac{1}{(c-1)!}}. \tag{4.73}$$

**Remark 4.14** From (4.60), when the SP upcrosses  $x$ , it next downcrosses  $x$  after a time  $a_x$  where  $E(a_x)$  is independent of  $x$ . By contrast, (4.73) implies when the SP downcrosses level  $x$ , it next upcrosses  $x$  after a time  $b_x$  where  $E(b_x)$  grows exponentially with increasing  $x$ .

The foregoing results for

$$d_x, a_x, b_x, E(d_x), E(a_x), E(b_x)$$

generalize analogous results for M/M/1 (Sect. 3.5.7).

### 4.9 M/M/c/c and Standard M/M/c Queues

The M/M/c/c queue is a special case of M/M/c/k, in which an upper limit  $k$  is placed on the number of customers allowed in the system at any time (see, e.g., Sect. 2.5, p. 76ff in [84]). Here, we develop a relationship between M/M/c/c and the standard  $M_\lambda/M_\mu/c$  queue. By a judicious choice of parameters for M/M/c/c, the pdf of the virtual wait in the two models have identical forms. However, the jump structure of the sample path of M/M/c/c is much simpler than that of the corresponding M/M/c model, for positive values of the virtual wait. This jump structure makes it much easier to derive the pdf of the virtual wait in M/M/c/c. The point of this exercise is to obtain the pdf of wait in the parameter-modified M/M/c/c queue, which can be derived in one line, without having to solve an integral equation (as in M/M/c), and the derived pdf is the same as in M/M/c. This relationship suggests a broader prospect. For a given complex model, can we identify a related model having the same solution form, that can be solved more easily?

The M/M/c/c queue is usually analyzed using a birth-death analysis. Here, we employ an LC approach. Consider an M/M/c/c queue where the service time for each customer that enters the system has exponential rate  $\mu - \frac{\lambda}{c} > 0$ . (We choose  $\lambda < c\mu$  because our related model is a standard  $M_\lambda/M_\mu/c$  queue in equilibrium.)

In M/M/c/c all *actual* waits are 0 – there is no waiting line. In a queue where blocking is possible, we shall define the virtual wait as the time that a potential arrival *would* wait to start service, if it were not blocked and cleared.

Thus the virtual wait is not 0 for every arrival. In M/M/c/c, customers that arrive when the *virtual* wait is positive, are blocked and cleared from the system. In both models, the virtual wait is positive if and only if all  $c$  servers are occupied.

For M/M/c/c, consider the ‘system point’ process  $\{W(t), M(t)\}_{t \geq 0}$ , where  $W(t)$  is the virtual wait and  $M(t) \in \{0, \dots, c - 1\}$  is the system configuration at time  $t$ .  $M(t)$  is the number of occupied servers at instant  $t^-$ , if there is an idle server at  $t^-$ . We denote the  $c$  discrete states by  $\{(0, 0), \dots, (0, c - 1)\}$ . Thus  $M(t) = n$  if  $n$  other servers are occupied when a customer joins the system and starts service,  $n = 0, \dots, c - 1$ . Denote the steady-state probability of  $(0, n)$  as  $P_n, n = 0, \dots, c - 1$ . Denote the *positive* virtual-wait states as  $\{(x, c - 1), x \in (0, \infty)\}$ .

#### 4.9.1 Sample Path of $\{W(t), M(t)\}_{t \geq 0}$

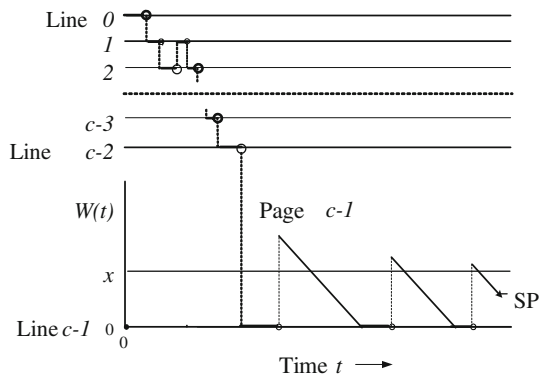
Consider a sample path of  $\{W(t), M(t)\}_{t \geq 0}$  (Fig. 4.5). Without loss of generality, assume the system starts empty. The SP is on line 0 at  $t = 0$ . As the system evolves, the SP moves among the lines until  $c - 1$  of the servers are occupied, just as in a standard  $M_\lambda/M_{\mu-\frac{\lambda}{c}}/c$  model. In Fig. 4.5 we situate line  $c - 1$  at level 0 of the page; this layout makes it easier to depict SP exchanges between line  $c - 1$  and the virtual-wait positive states.

Suppose a customer arrives when  $c - 1$  servers are occupied. The arrival joins the system and starts service in the one free server. All  $c$  servers are busy just after the arrival starts service. The configuration is  $c - 1$ , since  $c - 1$  other servers are occupied just after the arrival instant. Each of the  $c$  servers has service time  $\stackrel{dis}{=} \text{Exp}_{\mu-\frac{\lambda}{c}}$  once the arrival starts service, due to the memoryless property of exponential service times. The SP jumps to ordinate  $y \in (0, \infty)$  on the page, where  $y \stackrel{dis}{=} \text{Exp}_{c\mu-\lambda}$  which is distributed as the minimum of  $c$  i.i.d. exponential r.v.s each distributed with rate  $\mu - \frac{\lambda}{c}$ .

The SP descends at rate 1 (slope = -1), until it makes a continuous hit of level 0 from above. New arrivals are blocked and cleared, and have no effect on the sample path during this descent. Once the SP hits level 0, it continues its motion among the states  $(0, 0), \dots, (0, c - 1)$ , until it makes another jump out of state  $(0, c - 1)$  onto the page.

All upward jumps that end on the page start at level 0. Hence the jump structure for M/M/c/c is much simpler than that of the standard M/M/c queue, in which jumps that end on the page may start at any point in  $[0, \infty)$ .

**Fig. 4.5** Sample path of  $\{W(t)\}_{t \geq 0}$  in M/M/c/c queue. All jumps ending in  $(0, \infty)$  begin at level 0 (state  $(0, c - 1)$ ) with size  $\stackrel{dis}{=} \text{Exp}_{c\mu - \lambda}$



### 4.9.2 PDF of Virtual Wait

Denote the pdf of the virtual wait as  $f_{c-1}(x) \equiv f(x)$ ,  $x > 0$ . To derive the pdf of  $f(x)$ , fix level  $x > 0$ . The SP downcrossing rate of  $x$  is  $f(x)$ . Since all SP jumps ending on the page start from state  $(0, c - 1)$  at arrival instants, and all jumps sizes are  $\stackrel{dis}{=} \text{Exp}_{c\mu - \lambda}$ , the upcrossing rate of  $x$  is  $\lambda P_{c-1} e^{-(c\mu - \lambda)x}$ . Balancing SP rates out of and into set  $((x, \infty), c - 1)$  yields

$$f(x) = \lambda P_{c-1} e^{-(c\mu - \lambda)x}, x > 0. \tag{4.74}$$

**Remark 4.15** Formula (4.74) has precisely the same **form** as the steady-state pdf of wait in the standard  $M_\lambda/M_\mu/c$  queue given by (4.56), except that  $P_{c-1}$  has a **different value**. For the  $M_\lambda/M_{\mu - \lambda/c}/c/c$  queue, formula (4.74) is derived “instantly” from observing a sample path of the virtual wait. There is no need to solve an integral equation, as in M/M/c. In M/M/c/c, the pdf formula for  $f(x)$  is inherently a model equation. This is the main relationship between the two models we discuss here. The result for  $M_\lambda/M_{\mu - \lambda/c}/c/c$  allows us to write the **form** of the pdf of wait in  $M_\lambda/M_\mu/c$  immediately.

### 4.9.3 Non-blocking States

The rate-balance equations for the non-blocking states  $(0, 0), \dots, (0, c - 1)$  are the same as in (4.52) for  $M_\lambda/M_\mu/c$ , with  $\mu - \frac{\lambda}{c}$  substituted for  $\mu$ . Thus in M/M/c/c

$$P_n = \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^n \frac{1}{n!} P_0, n = 0, \dots, c - 1,$$

so that

$$P_{c-1} = \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^{c-1} \frac{1}{(c-1)!} P_0.$$

The normalizing condition is

$$\left( \sum_{n=0}^{c-1} \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^n \frac{1}{n!} \right) P_0 + \int_{x=0}^{\infty} f(x) dx = 1.$$

Applying (4.74) gives

$$\begin{aligned} & \left( \sum_{n=0}^{c-1} \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^n \frac{1}{n!} \right) P_0 \\ & + \lambda \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^{c-1} \frac{1}{(c-1)!} P_0 \int_{x=0}^{\infty} e^{-(c\mu - \lambda)x} dx = 1, \\ P_0 &= \frac{1}{\sum_{n=0}^{c-1} \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^n \frac{1}{n!} + \lambda \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^{c-1} \frac{1}{(c-1)!} \frac{1}{c\mu - \lambda}} \\ &= \frac{1}{\sum_{n=0}^c \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^n \frac{1}{n!}} \cdot \checkmark \end{aligned}$$

### 4.9.4 Blocking Time $T_B$

Let  $T_B$  denote the time from the instant the system gets blocked (all  $c$  servers occupied) until the first instant that it becomes unblocked thereafter (at which  $c - 1$  servers are occupied). We call  $T_B$  the blocking time.

The pdf of the virtual wait in  $M_\lambda/M_{\mu-\lambda}/c/c$  is the same as the pdf of  $S$  (inter start-of-service depart time) when an arrival “sees” state  $(0, c - 1)$ . Also,  $\mathcal{S} = T_B$ .

Then  $\overset{dis}{E}(T_B) = \text{Exp}_{c\mu - \lambda} = 1/(c\mu - \lambda)$ . Let  $P_c$  denote the proportion of time the system is blocked. Then

$$P_c = \int_{x=0}^{\infty} f(x) dx = \lambda P_{c-1} \int_{x=0}^{\infty} e^{-(c\mu - \lambda)x} dx$$



$$\begin{aligned}
&= \lambda \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^{c-1} \frac{1}{(c-1)!} P_0 \int_{x=0}^{\infty} e^{-(c\mu-\lambda)x} dx \\
&= \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^c \frac{1}{c!} P_0 \\
&= \frac{\left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^c \frac{1}{c!}}{\sum_{n=0}^c \left( \frac{\lambda}{\mu - \frac{\lambda}{c}} \right)^n \frac{1}{n!}}.
\end{aligned}$$

$P_c$  is the probability that a right-truncated Poisson variate (truncated at  $c$ ), has value  $c$ . It is the classical *Erlang-B loss formula* for the  $M_\lambda/M_{\mu-\frac{\lambda}{c}}/c/c$  queue (see, e.g., p. 82, Sect. 2.6 in [84]).

Note that the blocking time is a  $[c-1, c]$  busy period, denoted by  $\mathcal{B}_{c-1,c}$ , so that  $T_B \stackrel{dis}{=} \mathcal{B}_{c-1,c}$ . From Remark 4.8.4,  $E(\mathcal{B}_{c-1,c}) = \frac{1}{c\mu-\lambda}$ .

**Remark 4.16** Suppose that in the M/M/c/c model the servers were numbered  $1, \dots, c$ . Let the service rates assigned to arrivals depend on which server is occupied, say rates  $\nu_i, i = 1, \dots, c$ . Assume  $\sum_{i=1}^c \nu_i = c\mu - \lambda > 0$ , where  $\mu, \lambda$  are the parameters of a stable  $M_\lambda/M_\mu/c$  queue. Then the distribution of  $T_B$  would be the same as in (4.74). So this specialized M/M/c/c model can also be used as a “companion” model to obtain the pdf of wait in the  $M_\lambda/M_\mu/c$  queue.

### 4.9.5 Discussion

We can derive formula (4.74) for  $f(x)$  geometrically as in Sect. 4.8.4. Let  $F(x), x \geq 0$ , be the cdf of the virtual wait. We get

$$\begin{aligned}
\frac{d}{dx} \ln(1 - F(x)) &= \frac{-1}{E(\mathcal{B}_{c-1,c})} = -(c\mu - \lambda), \\
F(x) &= 1 - (1 - F(0))e^{-(c\mu-\lambda)x}, \quad x \geq 0, \\
f(x) &= (c\mu - \lambda)(1 - F(0))e^{-(c\mu-\lambda)x}. \quad (4.75)
\end{aligned}$$

Comparing (4.74) and (4.75) shows that

$$\lambda P_{c-1} = (c\mu - \lambda)(1 - F(0)) = (c\mu - \lambda) P_c, \quad (4.76)$$

where  $P_c$  is the probability of  $c$  units in the system.

In M/M/c/c an arrival enters the system iff the virtual wait is 0. Thus  $F(0) = P(\text{an arrival enters the system})$ . Hence  $(1 - F(0)) = P(\text{an arrival is blocked and cleared}) = P_c$ . Equation (4.76) is precisely the balance equation that would appear in a birth-death analysis of the system.

### 4.10 M/M/c in Which Zero-Wait Customers Get Special Service

Consider an M/M/c ( $c \geq 2$ ) queue with arrival rate  $\lambda$ , in which zero-wait customers get service rate  $\mu_0$ , and positive-wait customers get service rate  $\mu_1 (\neq \mu_0)$ . Thus, the assigned service rate is *state-dependent*. We derive below the steady-state pdf of wait, distribution of the number-in-system, and related model characteristics.

Denote the state of the system as  $\{W(t), \mathbf{M}(t)\}_{t \geq 0}$ , where  $W(t) \geq 0$  is the virtual wait and  $M(t)$  is the system configuration. Thus

$$\mathbf{M}(t) = (m_0, m_1), 0 \leq m_0 + m_1 \leq c - 1,$$

where  $m_j$  is the number of occupied servers operating at rate  $\mu_j, j = 0, 1$ . In the notation of Sect. 4.4, integer  $J = 1$ . The number of zero-wait states is the total number of non-negative integer solutions for  $m_0, m_1$  in the  $c$  equations

$$m_0 + m_1 = k, k = 0, \dots, c - 1,$$

which is, since  $J = 1$ ,

$$\begin{aligned} \sum_{k=0}^{c-1} \binom{J+k}{J} &= \binom{J+c}{J+1} = \binom{c+1}{2} \\ &= \frac{c(c+1)}{2} = 1 + 2 + \dots + c. \end{aligned}$$

From Sect. 4.4.2,  $\mathbf{M}_0 = \left\{ (0, \mathbf{m}) \mid 0 \leq \sum_{j=0}^J m_j \leq c - 2 \right\}$ , which contains  $\frac{(c-1)c}{2}$  configurations. Set  $\mathbf{M}_b = \left\{ \mathbf{m} \mid \sum_{j=0}^J m_j = c - 1 \right\}$  comprises the discrete *boundary* states, and contains  $\binom{J+c-1}{J} = \binom{c}{1} = c$  configurations. (Note that  $\mathbf{M}_b = \mathbf{M}_1$ .)

**Zero-wait Probabilities** Let  $P_{m_0 m_1}$  denote the steady-state probability that an arrival “sees”  $m_j$  rate- $\mu_j$  customers in service,  $j = 0, 1$ , and waits

zero before starting service.  $P_{m_0 m_1}$  is the steady-state probability of state  $(0, (m_0, m_1))$ .

There are  $c$  positive-wait *pages (sheets)*, one for each configuration in  $\mathbf{M}_b$ , where

$$\mathbf{M}_b = \{(c-1, 0), (c-2, 1), \dots, (1, c-2), (0, c-1)\}.$$

**Positive-wait PDFs** Let  $f_m(x)$ ,  $x > 0$ , denote the steady-state pdf of the virtual wait when the occupancies of the *other*  $c-1$  servers will be  $\mathbf{m} \in \mathbf{M}_b$  at start of service ('look-ahead' property of virtual wait).

### 4.10.1 Equations for Probabilities of Zero-Wait States

The  $\frac{c(c+1)}{2}$  zero-wait states, having configurations in  $\mathbf{M}_0 \cup \mathbf{M}_b$ , viz.,

$$(0, (m_0, m_1)), 0 \leq m_0 + m_1 \leq c-1,$$

give rise to  $\frac{c(c+1)}{2}$  linear equations for their probabilities, using the principle *rate out = rate in*, as in (4.77)–(4.79) below.

First consider states  $(0, \mathbf{m})$ ,  $\mathbf{m} \in \mathbf{M}_0$ . For  $m_0 = m_1 = 0$ , (empty system) there is one equation:

$$\lambda P_{00} = \mu_0 P_{10} + \mu_1 P_{01}. \quad (4.77)$$

For states  $(0, (m_0, m_1))$ ,  $1 \leq m_0 + m_1 \leq c-2$ , there are  $\frac{(c-1)c}{2} - 1$  equations, each of the form

$$\begin{aligned} (\lambda + m_0 \mu_0 + m_1 \mu_1) P_{m_0 m_1} = & \lambda P_{(m_0-1)m_1} \\ & + (m_0 + 1) \mu_0 P_{(m_0+1)m_1} \\ & + (m_1 + 1) \mu_1 P_{m_0(m_1+1)}. \end{aligned} \quad (4.78)$$

For states  $(0, (m_0, m_1)) \in \mathbf{M}_b$ , there are  $c$  equations, each of the form

$$(\lambda + m_0 \mu_0 + m_1 \mu_1) P_{m_0 m_1} = \lambda P_{(m_0-1)m_1} + f_{m_0 m_1}(0). \quad (4.79)$$

In (4.79) the term  $f_{m_0 m_1}(0)$  ( $= f_{m_0 m_1}(0^+)$ ) is the rate at which the SP enters border state  $(0, (m_0, m_1))$  due to left continuous hits of level 0 from above on page  $m_0 m_1$ .

### 4.10.2 Equations for PDF of Positive-Wait States

There are  $c$  Volterra integral equations for the positive-wait states. Consider composite state  $((x, \infty), \mathbf{m})$ ,  $x > 0$ , on page  $\mathbf{m} \in \mathbf{M}_b$ . For positive-wait states  $(y, m_0 m_1)$ ,  $y > 0$ ,  $m_0 + m_1 = c - 1$ . We first specify the SP exit and entrance rates of the pertinent composite states in the state space. Then we will write the equations.

**Rate Out of  $((x, \infty), m_0 m_1)$**

Because  $(m_0, m_1)$  is a configuration,  $m_0 + m_1 = c - 1$ . The SP rate out of  $((x, \infty), m_0 m_1)$  is

$$f_{m_0 m_1}(x) + \lambda \frac{m_0 \mu_0}{m_0 \mu_0 + (m_1 + 1) \mu_1} \int_{y=x}^{\infty} f_{m_0 m_1}(y) dy. \quad (4.80)$$

#### Explanation of Terms in (4.80)

The first term  $f_{m_0 m_1}(x)$  is the SP downcrossing rate of level  $x$  on page  $m_0 m_1$ . The second term

$$\lambda \frac{m_0 \mu_0}{m_0 \mu_0 + (m_1 + 1) \mu_1} \int_{y=x}^{\infty} f_{m_0 m_1}(y) dy$$

is the rate of arrivals when the state is  $(y, m_0 m_1)$ ,  $y > x$  (being assigned service rate  $\mu_1$  thereby adding one rate- $\mu_1$  occupied server *upon start of service*); and a rate- $\mu_0$  service completes first thereafter. At the arrival instant the SP jumps to level  $y + \text{Exp}_{m_0 \mu_0 + (m_1 + 1) \mu_1}$  on page  $(m_0 - 1, m_1 + 1)$  (i.e., page  $(m_0 - 1, c - m_0)$ ). If  $m_0 = 0$ , the SP would be on page  $(0, c - 1)$ . The only exit from the page would be via a downcrossing of level 0. All arrivals would be assigned service rate  $\mu_1$  and cause the SP to jump upward but remain on page  $(0, c - 1)$ ; the second term in (4.80) would equal 0 if  $m_0 = 0$ .

**Rate into  $((x, \infty), m_0 m_1)$**

The SP rate into  $((x, \infty), m_0 m_1)$  is

$$\begin{aligned} & \lambda \frac{(m_0 + 1) \mu_0}{(m_0 + 1) \mu_0 + m_1 \mu_1} e^{-((m_0 + 1) \mu_0 + m_1 \mu_1) x} P_{m_0 m_1} \\ & + \lambda \frac{(m_1 + 1) \mu_1}{m_0 \mu_0 + (m_1 + 1) \mu_1} e^{-(m_0 \mu_0 + (m_1 + 1) \mu_1) x} P_{m_0 - 1, m_1 + 1} \\ & + \lambda \frac{(m_0 + 1) \mu_0}{(m_0 + 1) \mu_0 + m_1 \mu_1} \int_{y=x}^{\infty} f_{m_0 + 1, m_1 - 1}(y) dy \\ & + \lambda \frac{(m_0 + 1) \mu_0}{(m_0 + 1) \mu_0 + m_1 \mu_1} \int_{y=0}^x e^{-((m_0 + 1) \mu_0 + m_1 \mu_1)(x - y)} f_{m_0 + 1, m_1 - 1}(y) dy \end{aligned}$$

$$+ \lambda \frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1} \int_{y=0}^x e^{-(m_0\mu_0 + (m_1 + 1)\mu_1)(x-y)} f_{m_0 m_1}(y) dy. \quad (4.81)$$

where we have inserted a comma in subscripts like  $m_0 - 1, m_1 + 1$ , for clarity.

### Explanation of Terms in (4.81)

The term

$$\lambda \frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1} e^{-((m_0 + 1)\mu_0 + m_1\mu_1)x} P_{m_0 m_1}$$

is the rate at which the SP jumps at arrival instants from level 0 on page  $m_0 m_1$  into  $((x, \infty), m_0 m_1)$ . At arrival instants customers are assigned service rate  $\mu_0$  (wait = 0), resulting in  $(m_0 + 1)$  rate- $\mu_0$  and  $m_1$  rate- $\mu_1$  customers in service. A rate- $\mu_0$  service finishes first with probability

$$\frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1},$$

in which case the SP jumps to page  $m_0 m_1$ . SP jumps from level 0 over level  $x$  have probability  $e^{-((m_0 + 1)\mu_0 + m_1\mu_1)x}$  since  $\mathcal{S} \stackrel{dis}{=} \text{Exp}_{(m_0 + 1)\mu_0 + m_1\mu_1}$ .

The term

$$\lambda \frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1} e^{-(m_0\mu_0 + (m_1 + 1)\mu_1)x} P_{m_0 - 1, m_1 + 1}$$

is the rate at which the SP jumps at arrival instants, from level 0 on page  $(m_0 - 1, m_1 + 1)$  into  $((x, \infty), m_0 m_1)$ . The arriving customer is assigned service rate  $\mu_0$  (wait = 0), resulting in  $m_0$  rate- $\mu_0$  and  $(m_1 + 1)$  rate- $\mu_1$  customers in service. If a rate- $\mu_1$  service finishes first thereafter, the SP jumps to page  $m_0 m_1$ ; the probability is

$$\frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1}.$$

SP jumps from level 0 upcross level  $x$  with probability  $e^{-(m_0\mu_0 + (m_1 + 1)\mu_1)x}$  since the inter-start-of-service depart time  $\mathcal{S} \stackrel{dis}{=} \text{Exp}_{m_0\mu_0 + (m_1 + 1)\mu_1}$ .

The term

$$\lambda \frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1} \int_{y=x}^{\infty} f_{m_0 + 1, m_1 - 1}(y) dy$$

is the rate at which the SP jumps at arrival instants, out of  $(x, \infty)$  on page  $(m_0 + 1, m_1 - 1)$  into  $((x, \infty), m_0 m_1)$ . The arriving customer is assigned service rate  $\mu_1$  (wait  $> 0$ ) resulting in  $(m_0 + 1)$  rate- $\mu_0$  and  $m_1$  rate- $\mu_1$  customers in service just after the start of service of the arrival. If a rate- $\mu_0$  service finishes first, the SP jumps to page  $m_0 m_1$ ; this has probability

$$\frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1}.$$

A jump  $\mathcal{S}$  of any size will cause such a jump to enter  $((x, \infty), m_0 m_1)$  since the start of the jump is already above level  $x$ .

The term

$$\lambda \frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1} \int_{y=0}^x e^{-((m_0+1)\mu_0+m_1\mu_0)(x-y)} f_{m_0+1, m_1-1}(y) dy$$

is the rate at which the SP jumps upward at arrivals, out of

$$((0, x), (m_0 + 1, m_1 - 1)) \text{ into } ((x, \infty), m_0 m_1).$$

That is, the SP makes a  $(m_0 + 1, m_1 - 1) \rightarrow (m_0 m_1)$  upcrossing of level  $x$ . An arrival is assigned service rate  $\mu_1$  (wait  $> 0$ ). Just after the arrival starts service there are  $m_0 + 1$  rate- $\mu_0$  and  $m_1$  rate- $\mu_1$  customers in service. The probability that a rate- $\mu_0$  service finishes first is

$$\frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1},$$

causing the SP to jump to page  $m_0 m_1$ . Starting at level  $y < x$  the SP will upcross level  $x$  if  $\mathcal{S} > x - y$ ; since  $\mathcal{S} \stackrel{\text{dis}}{=} \text{Exp}_{(m_0+1)\mu_0+m_1\mu_1}$ , this event has probability  $e^{-((m_0+1)\mu_0+m_1\mu_1)(x-y)}$ .

The term

$$\lambda \frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1} \int_{y=0}^x e^{-(m_0\mu_0+(m_1+1)\mu_1)(x-y)} f_{m_0 m_1}(y) dy$$

is the rate at which the SP jumps at arrival instants from  $((0, x), m_0 m_1)$  upward into  $((x, \infty), m_0 m_1)$ , i.e., it upcrosses level  $x$  on page  $m_0 m_1$ . Arrivals are assigned service rate  $\mu_1$  (wait  $> 0$ ). Just after the arrival starts service there are  $m_0$  rate- $\mu_0$  and  $(m_1 + 1)$  rate- $\mu_1$  customers in service. A rate- $\mu_1$  service ends first with probability

$$\frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1}.$$

causing the SP to jump to page  $m_0m_1$ . If the SP starts at level  $y$  it will upcross level  $x$  provided  $\mathcal{S} > x - y$ ; since  $\mathcal{S} \stackrel{dis}{=} \text{Exp}_{m_0\mu_0 + (m_1 + 1)\mu_1}$ , this event has probability  $e^{-(m_0\mu_0 + (m_1 + 1)\mu_1)(x - y)}$ .

### Writing Equations for Positive-Wait States

The model equation for the positive-wait states on page  $m_0m_1$  is written by using the principle of rate balance with respect to set  $((x, \infty), m_0m_1)$ , *exit rate = entrance rate*. Equating exit rate (4.80) and entrance rate (4.81) gives

$$\begin{aligned} & f_{m_0m_1}(x) + \lambda \frac{m_0\mu_0}{m_0\mu_0 + (m_1 + 1)\mu_1} \int_{y=x}^{\infty} f_{m_0m_1}(y) dy \\ &= \lambda \frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1} e^{-((m_0 + 1)\mu_0 + m_1\mu_1)x} P_{m_0m_1} \\ &+ \lambda \frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1} e^{-(m_0\mu_0 + (m_1 + 1)\mu_1)x} P_{m_0 - 1, m_1 + 1} \\ &+ \lambda \frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1} \int_{y=x}^{\infty} f_{m_0 + 1, m_1 - 1}(y) dy \\ &+ \lambda \frac{(m_0 + 1)\mu_0}{(m_0 + 1)\mu_0 + m_1\mu_1} \int_{y=0}^x e^{-((m_0 + 1)\mu_0 + m_1\mu_1)(x - y)} f_{m_0 + 1, m_1 - 1}(y) dy \\ &+ \lambda \frac{(m_1 + 1)\mu_1}{m_0\mu_0 + (m_1 + 1)\mu_1} \int_{y=0}^x e^{-(m_0\mu_0 + (m_1 + 1)\mu_1)(x - y)} f_{m_0m_1}(y) dy. \end{aligned} \tag{4.82}$$

### Equation for “Cover”

The total probability of a zero wait is

$$P_0 = \sum_{m \in M_0 \cup M_b} P_m = \sum_{0 \leq m_0 + m_1 \leq c - 1} P_{m_0m_1}. \tag{4.83}$$

The total pdf of wait is

$$f(x) = \sum_{m \in M_1} f_m(x) = \sum_{m_0 + m_1 = c - 1} f_{m_0m_1}(x), \quad x > 0. \tag{4.84}$$

Let  $x > 0$  be fixed. The total SP downcrossing rate of  $x$  is  $f(x)$ . The total SP upcrossing rate of  $x$  due to jumps starting from level 0 at arrival instants, is

$$\lambda \sum_{m_0+m_1=c-1} e^{-((m_0+1)\mu_0+m_1\mu_1)x} P_{m_0m_1}.$$

The total SP upcrossing rate of  $x$  due to jumps starting from levels  $y \in (0, x)$  at arrival instants, is

$$\lambda \sum_{m_0+m_1=c-1} \int_{y=0}^x e^{-(m_0\mu_0+(m_1+1)\mu_1)(x-y)} f_{m_0m_1}(y) dy.$$

Rate balance across level  $x$  gives the model equation for the cover,

$$\begin{aligned} f(x) = & \lambda \sum_{m_0+m_1=c-1} e^{-((m_0+1)\mu_0+m_1\mu_1)x} P_{m_0m_1} \\ & + \lambda \sum_{m_0+m_1=c-1} \int_{y=0}^x e^{-(m_0\mu_0+(m_1+1)\mu_1)(x-y)} f_{m_0m_1}(y) dy. \end{aligned} \quad (4.85)$$

### Normalizing Condition

The normalizing condition  $P_0 + \int_{x=0}^{\infty} f(x) dx = 1$  can be expressed as

$$\sum_{0 \leq m_0+m_1 \leq c-1} P_{m_0m_1} + \sum_{m_0+m_1=c-1} \int_{x=0}^{\infty} f_{m_0m_1}(x) dx = 1. \quad (4.86)$$

### 4.10.3 Solution of Model Equations

In Sect. 4.11 below, we formulate and solve the foregoing M/M/2 model with zero-wait customers receiving exceptional service, whose solution illustrates relevant SPLC ideas and related insights. A more general solution procedure of a two-server M/M/2 queue where service time depends on waiting time in a *general manner* is detailed in Chap. 4 of [11].

## 4.11 M/M/2: Zero-Waits Get Special Service

**M/M/2/** $(\mu_0, \mu_1), (0, (0, \infty))$

To fix ideas and clarify the system dynamics of M/M/c with special service for zero-wait customers, we formulate the model with  $c = 2$  servers. We discuss



the solution for the zero-wait probabilities and the positive-wait pdfs. We denote the model by  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$ . This notation indicates that 0-wait arrivals get service rate  $\mu_0$  and  $(0, \infty)$ -wait arrivals get service rate  $\mu_1$ ; diagrammatically,  $\mu_0 \leftrightarrow 0$ -wait,  $\mu_1 \leftrightarrow (0, \infty)$ -wait.

There are only three zero-wait states in  $M_0 \cup M_b$  (compare Sect. 4.10),

$$\{(0, m_0 m_1)\} = \{(0, 00), (0, 10), (0, 01)\}.$$

Denote the steady-state probabilities of the zero-wait states by  $P_{00}, P_{10}, P_{01}$  respectively.

For example, state  $(0, 10)$  indicates that an arrival would wait 0 and would “see” a rate- $\mu_0$  customer being served by the *other* server. The arrival would be assigned rate  $\mu_0$  since it waits 0. There would then be two rate- $\mu_0$  customers in service. The inter start-of-service depart time  $\mathcal{S}$  would be  $\underset{dis}{=} \text{Exp}_{2\mu_0}$ .

There are only two zero-wait states such that  $m_0 + m_1 = 1$  (both border states). Denote the pdfs of the positive-wait states  $(x, 10), (x, 01)$ , by  $f_{10}(x), f_{01}(x), x > 0$ , respectively. A would-be arrival that finds the state  $(x, 10), x > 0$ , for example, would wait  $x$  before service, and be assigned service rate  $\mu_1$  (wait  $> 0$ ). Just after its start of service, it would have a rate- $\mu_0$  customer as neighbor in the other server. Inter start-of-service depart time  $\mathcal{S} \underset{dis}{=} \text{Exp}_{\mu_0 + \mu_1}$ . The rate- $\mu_1$  customer would finish service first with probability  $\frac{\mu_1}{\mu_0 + \mu_1}$ , leaving the rate- $\mu_0$  customer in service ( $m_0 m_1 = 10$ ). The rate- $\mu_0$  customer would finish service first with probability  $\frac{\mu_0}{\mu_0 + \mu_1}$ , leaving the rate- $\mu_1$  customer in service ( $m_0 m_1 = 01$ ).

If an arrival “sees” state  $(x, 01), x > 0, \mathcal{S}$  would be  $\underset{dis}{=} \text{Exp}_{2\mu_1}$ . The first customer to complete service would have rate  $\mu_1$  with certainty. The customer remaining in service just after that service completion would have service rate  $\mu_1$  ( $m_0 m_1 = 01$ ).

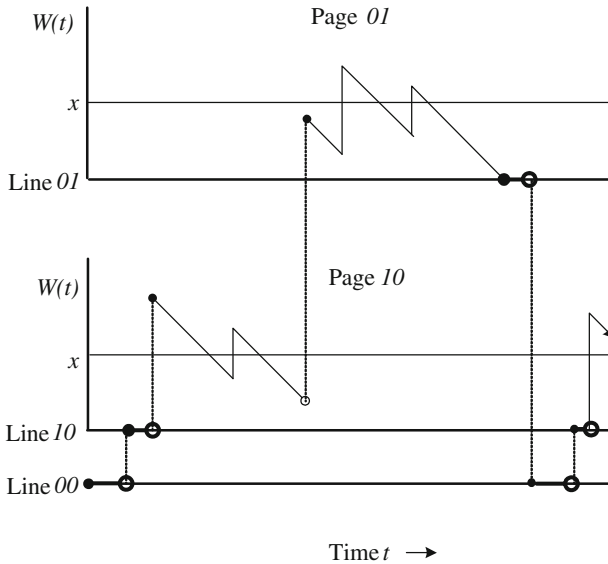
A sample-path diagram of the virtual wait process  $\{W(t)\}_{t \geq 0}$ , has three lines and two pages (Fig. 4.6).

The total (marginal) probability of a zero wait is

$$P_0 = P_{00} + P_{10} + P_{01}.$$

The total pdf of wait is

$$f(x) = f_{10}(x) + f_{01}(x), x > 0.$$



**Fig. 4.6** Sample path of virtual wait in  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$ . Lines for states  $(0, 10), (0, 01)$  are at level 0 of corresponding pages. Line for state  $(0, 00)$  is isolated. The SP can enter state  $(0, 01)$  only by downcrossing level 0 on page 01 (See Fig. 4.10.)

### 4.11.1 Model Equations

#### Zero-Wait States

Applying *SP exit rate = SP entrance rate* for the zero-wait states  $(0, 00), (0, 10), (0, 01)$  gives, respectively,

$$\begin{aligned} \lambda P_{00} &= \mu_0 P_{10} + \mu_1 P_{01}, \\ (\lambda + \mu_0) P_{10} &= \lambda P_{00} + f_{10}(0), \\ (\lambda + \mu_1) P_{01} &= f_{01}(0). \end{aligned} \tag{4.87}$$

In (4.87), the terms  $f_{10}(0), f_{01}(0)$  (same as  $f_{10}(0^+), f_{01}(0^+)$ ) are the rates at which the SP hits level 0 from above on pages 10 and 01 respectively. Immediately following such hits, the SP moves on lines 10 and 01 respectively.

#### Positive-Wait States

Applying *SP exit rate = SP entrance rate* for  $((x, \infty), 10)$  (on page 10) yields the integral equation

$$\begin{aligned}
& f_{10}(x) + \lambda \int_{y=x}^{\infty} \frac{\mu_0}{\mu_0 + \mu_1} f_{10}(y) dy \\
&= \lambda P_{10} e^{-2\mu_0 x} + \lambda \frac{\mu_1}{\mu_0 + \mu_1} P_{01} e^{-(\mu_0 + \mu_1)x} \\
&\quad + \lambda \frac{\mu_1}{\mu_0 + \mu_1} \int_{y=0}^x e^{-(\mu_0 + \mu_1)(x-y)} f_{10}(y) dy, \quad x > 0. \quad (4.88)
\end{aligned}$$

When formulating equation (4.88), note that the SP cannot jump directly from a positive-wait state on page 01 into set  $((x, \infty), 10)$ . An arrival that “sees” state  $(y, 01)$ ,  $y > 0$ , will be assigned rate  $\mu_1$  and start service after a wait  $y$ ; its neighbor in the other server will also have service rate  $\mu_1$  (because  $m_0 m_1 = 01$ ). The random variable  $\mathcal{S}$  will be distributed as  $\text{Exp}_{2\mu_1}$ , and the remaining customer in service just after the first departure thereafter, will have rate  $\mu_1$ . At the arrival instant, the SP will start a jump at level  $y$  on page 01, which ends at level  $y + \mathcal{S} = y + \text{Exp}_{2\mu_1}$ , also on page 01. The configuration remains 01 just after the arrival. The only exit route from page 01 is via a downcrossing of level 0 (continuous hit of 0 from above—see Fig. 4.6).

Now we balance the SP exit and entrance rates for  $((x, \infty), 01)$  (page 01), giving integral equation

$$\begin{aligned}
f_{01}(x) &= \lambda \frac{\mu_0}{\mu_0 + \mu_1} P_{01} e^{-(\mu_0 + \mu_1)x} \\
&\quad + \lambda \int_{y=0}^x e^{-2\mu_1(x-y)} f_{01}(y) dy \\
&\quad + \lambda \frac{\mu_0}{\mu_0 + \mu_1} \int_{y=0}^x e^{-(\mu_0 + \mu_1)(x-y)} f_{10}(y) dy \\
&\quad + \lambda \frac{\mu_0}{\mu_0 + \mu_1} \int_{y=x}^{\infty} f_{10}(y) dy. \quad (4.89)
\end{aligned}$$

When formulating (4.89), note that the SP can exit  $((x, \infty), 01)$  only by downcrossing level  $x$ . Also, the SP cannot enter  $((x, \infty), 01)$  from state  $(0, 10)$  at arrivals, since all jumps that start from line 10 (corresponding to state  $(0, 10)$ ) must end on page 10, at an ordinate  $= \text{Exp}_{2\mu_0}$ .

The equation for the total pdf is

$$f(x) = f_{10}(x) + f_{01}(x),$$

as viewed from the “cover”, the result of projecting sample-path segments on pages 10 and 01 onto a single sheet. An integral equation for  $f(x)$  is obtained by balancing the SP *total* down- and upcrossing rates of level  $x > 0$ . This is equivalent to equating the exit and entrance rates for the state-space set

$$((x, \infty), 10) \cup ((x, \infty), 01).$$

The resulting equation is

$$\begin{aligned} f(x) &= \lambda P_{10} e^{-2\mu_0 x} + \lambda P_{01} e^{-(\mu_0 + \mu_1)x} \\ &+ \lambda \int_{y=0}^x e^{-(\mu_0 + \mu_1)(x-y)} f_{10}(y) dy \\ &+ \lambda \int_{y=0}^x e^{-2\mu_1(x-y)} f_{01}(y) dy, \quad x > 0. \end{aligned} \quad (4.90)$$

Equation (4.90) can also be derived by summing the corresponding sides of (4.88) and (4.89). However, it is intuitive and instructive to interpret equation (4.90) as total SP rate-balance across level  $x > 0$ .

The normalizing condition is

$$P_{00} + P_{10} + P_{01} + \int_{x=0}^{\infty} f_{10}(x) dx + \int_{x=0}^{\infty} f_{01}(x) dx = 1,$$

or

$$P_0 + \int_{x=0}^{\infty} f(x) dx = 1. \quad (4.91)$$

### 4.11.2 Solution of Equations

Equation (4.88) is an integral equation in  $f_{10}(x)$ , which is not confounded by the presence of  $f_{01}(x)$ ; so we utilize it to obtain the functional form of  $f_{10}(x)$ . Applying differential operator  $\langle D \rangle \langle D + \mu_0 + \mu_1 \rangle$  to both sides of (4.88) leads to the second order differential equation

$$\begin{aligned} f_{10}''(x) + (\mu_0 + \mu_1 - \lambda) f_{10}'(x) - \lambda \mu_0 f_{10}(x) \\ = 2\lambda \mu_0 (\mu_0 - \mu_1) P_{10} e^{-2\mu_0 x}, \quad x > 0. \end{aligned} \quad (4.92)$$

with solution

$$f_{10}(x) = C_{10}e^{ax} + C_{10}^1e^{bx} + \lambda K_{10}P_{10}e^{-2\mu_0x}, \quad x > 0,$$

where  $a$  and  $b$  are the roots of the auxiliary quadratic equation, namely

$$\begin{aligned} a &= \frac{1}{2} \left( \lambda - \mu_0 - \mu_1 - \sqrt{\lambda^2 + 2\lambda\mu_0 - 2\lambda\mu_1 + \mu_0^2 + 2\mu_0\mu_1 + \mu_1^2} \right) < 0, \\ b &= \frac{1}{2} \left( \lambda - \mu_0 - \mu_1 + \sqrt{\lambda^2 + 2\lambda\mu_0 - 2\lambda\mu_1 + \mu_0^2 + 2\mu_0\mu_1 + \mu_1^2} \right) > 0, \\ K_{10} &= \frac{2(\mu_0 - \mu_1)}{\lambda + 2\mu_0 - 2\mu_1}, \end{aligned}$$

and  $C_{10}$ ,  $C_{10}^1$  are constants of integration. A necessary condition for system stability is  $\lim_{x \rightarrow \infty} f_{10}(x) = 0$ , which implies  $C_{10}^1 = 0$  (since  $b > 0$ ). Thus the functional form of  $f_{10}(x)$  is

$$f_{10}(x) = C_{10}e^{ax} + \lambda K_{10}P_{10}e^{-2\mu_0x}, \quad x > 0, \quad (4.93)$$

where  $C_{10}$  is a constant to be determined.

The term  $K_{10}$  will be undefined if  $\lambda + 2\mu_0 - 2\mu_1 = 0$ . If  $\lambda + 2\mu_0 - 2\mu_1 \neq 0$  and  $\mu_0 - \mu_1 \neq 0$ , then  $K_{10}$  may be positive or negative. If  $\mu_0 - \mu_1 = 0$  the model reduces to a standard M/M/c queue with  $c = 2$  (Sect. 4.8); the computed distribution of wait should then match that of a standard M/M/2 queue. (We will utilize this property later as a mild check on the correctness of the present solution.)

We obtain the functional form of  $f_{01}(x)$  by substituting the expression for  $f_{10}(x)$  (4.93) into (4.90). Since

$$f_{01}(x) = f(x) - f_{10}(x),$$

this substitution gives the integral equation

$$\begin{aligned} f_{01}(x) &= \lambda(1 - K_{10})P_{10}e^{-2\mu_0x} + \lambda P_{01}e^{-(\mu_0+\mu_1)x} - C_{10}e^{ax} \\ &\quad + \lambda \int_{y=0}^x e^{-(\mu_0+\mu_1)(x-y)} (C_{10}e^{ay} + \lambda K_{10}P_{10}e^{-2\mu_0y}) dy \\ &\quad + \lambda \int_{y=0}^x e^{-2\mu_1(x-y)} f_{01}(y) dy. \end{aligned} \quad (4.94)$$

The first integral term in (4.94) is

$$\begin{aligned} & \lambda \int_{y=0}^x e^{-(\mu_0+\mu_1)(x-y)} (C_{10}e^{ay} + \lambda K_{10}P_{10}e^{-2\mu_0y}) dy \\ &= \frac{\lambda C_{10}}{\mu_0 + \mu_1 + a} e^{ax} - \frac{\lambda^2 K_{10}P_{10}e^{-2\mu_0x}}{\mu_0 - \mu_1} \\ & \quad - \left( \frac{\lambda C_{10}}{\mu_0 + \mu_1 + a} - \frac{\lambda^2 K_{10}P_{10}}{\mu_0 - \mu_1} \right) e_1^{-(\mu_0+\mu_1)x}. \end{aligned}$$

Thus (4.94) is equivalent to the integral equation

$$\begin{aligned} f_{01}(x) &= H_{01}C_{10}e^{ax} + \lambda B_{01}P_{10}e^{-2\mu_0x} \\ & \quad + D_{01}e^{-(\mu_0+\mu_1)x} \\ & \quad + \lambda \int_{y=0}^x e^{-2\mu_1(x-y)} f_{01}(y) dy, \end{aligned} \quad (4.95)$$

where

$$\begin{aligned} H_{01} &= \frac{\lambda}{\mu_0 + \mu_1 + a} - 1, \\ B_{01} &= 1 - K_{10} - \frac{\lambda K_{10}}{\mu_0 - \mu_1}, \\ D_{01} &= \lambda P_{01} - \frac{\lambda C_{10}}{\mu_0 + \mu_1 + a} + \frac{\lambda^2 K_{10}P_{10}}{\mu_0 - \mu_1}. \end{aligned}$$

Applying the differential operator  $\langle D + 2\mu_1 \rangle$  to both sides of (4.95) yields the differential equation for  $f_{01}(x)$ ,

$$\begin{aligned} f'_{01}(x) + (2\mu_1 - \lambda)f_{01}(x) &= (2\mu_1 + a)H_{01}C_{10}e^{ax} \\ & \quad + 2\lambda(\mu_1 - \mu_0)B_{01}P_{10}e^{-2\mu_0x} \\ & \quad + (\mu_1 - \mu_0)D_{01}e^{-(\mu_1+\mu_0)x}. \end{aligned} \quad (4.96)$$

whose solution is

$$\begin{aligned} f_{01}(x) &= \frac{2\lambda(\mu_1 - \mu_0)}{2\mu_1 - \lambda - 2\mu_0} B_{01}P_{10}e^{-2\mu_0x} \\ & \quad + \frac{\mu_1 - \mu_0}{\mu_1 - \lambda - \mu_0} D_{01}e^{-(\mu_1+\mu_0)x} \end{aligned}$$

$$\begin{aligned}
 &+ \frac{2\mu_1 + a}{2\mu_1 - \lambda + a} H_{01} C_{10} e^{ax} \\
 &+ C_{01} e^{-(2\mu_1 - \lambda)x}, \tag{4.97}
 \end{aligned}$$

where  $C_{01}$  is a constant of integration to be determined (see Sect. 4.11.4).

### 4.11.3 Stability Condition

Consider the functional forms of  $f_{10}(x)$  and  $f_{01}(x)$  in (4.93) and (4.97). In the exponents, all the coefficients of  $x$  are negative except possibly the coefficient  $-(2\mu_1 - \lambda)$  in  $e^{-(2\mu_1 - \lambda)x}$  of (4.97). A necessary condition for stability is that

$$f_{10}(\infty) = f_{01}(\infty) = f(\infty) = 0;$$

implying  $-(2\mu_1 - \lambda) < 0$ , equivalent to  $\lambda < 2\mu_1$ . That is, the arrival rate must be less than the system departure rate when both servers are occupied by positive-wait customers, regardless how large  $x$  is. Thus, for stability, if the waiting time is large and customers are arriving, then the mean inter-arrival time should exceed the mean inter-departure time. This ensures that the waiting time will return to zero in a finite time.

### 4.11.4 Determination of Constants

A complete solution for the distribution of wait requires the values of five unknown constants

$$P_{00}, P_{10}, P_{01}, C_{10}, C_{01},$$

which we obtain from five independent equations.

In (4.93) letting  $x \downarrow 0$  to obtain  $f_{10}(0)$ , and referring to (4.87) gives

$$C_{10} + \lambda K_{10} P_{10} = (\lambda + \mu_0) P_{10} - \lambda P_{00}. \tag{4.98}$$

In (4.97) letting  $x \downarrow 0$  to obtain  $f_{01}(0)$  gives

$$\begin{aligned}
 f_{01}(0) &= \frac{2\lambda(\mu_1 - \mu_0)}{2(\mu_1 - \mu_0) - \lambda} B_{01} P_{10} \\
 &+ \frac{\mu_1 - \mu_0}{\mu_1 - \mu_0 - \lambda} D_{01} \\
 &+ \frac{2\mu_1 + a}{2\mu_1 + a - \lambda} H_{01} C_{10} + C_{01}. \tag{4.99}
 \end{aligned}$$

Substituting  $f_{01}(0)$  from (4.99) into (4.87) gives

$$\begin{aligned} C_{01} &= (\lambda + \mu_1)P_{01} - \frac{2\lambda(\mu_1 - \mu_0)}{2(\mu_1 - \mu_0) - \lambda} B_{01}P_{10} \\ &\quad - \frac{\mu_1 - \mu_0}{\mu_1 - \mu_0 - \lambda} D_{01} \\ &\quad - \frac{2\mu_1 + a}{2\mu_1 + a - \lambda} H_{01}C_{10}. \end{aligned} \quad (4.100)$$

We get another independent equation by substituting the functional form

$$f_{10}(x) = C_{10}e^{ax} + \lambda K_{10}P_{10}e^{-2\mu_0x}$$

into the integral equation (4.88) and equating the coefficients of corresponding exponential terms on both sides after evaluating the integral (different exponentials are linearly independent—see, e.g., Sect. 3.3, pp. 99ff and p. 205 in [10]). The coefficient of  $e^{-(\mu_0+\mu_1)x}$  on the right side of (4.88) must be 0. This yields the linear equation

$$\frac{\lambda\mu_1}{\mu_0 + \mu_1}P_{01} - \frac{1}{\mu_0 + \mu_1 + a} - \frac{\lambda K_{10}}{\mu_1 - \mu_0}P_{10} = 0. \quad (4.101)$$

The normalizing condition is

$$\begin{aligned} 1 &= P_{00} + P_{10} + P_{01} + \frac{C_{10}}{(-a)} + \frac{\lambda K_{10}P_{10}}{2\mu_0} \\ &\quad + \frac{\lambda(\mu_1 - \mu_0)}{\mu_0(2(\mu_1 - \mu_0) - \lambda)} B_{01}P_{10} + \frac{\mu_1 - \mu_0}{(\mu_1 + \mu_0)(\mu_1 - \mu_0 - \lambda)} D_{01} \\ &\quad + \frac{2\mu_1 + a}{(-a)(2\mu_1 + a - \lambda)} H_{01}C_{10} + \frac{1}{2\mu_1 - \lambda} C_{01}. \end{aligned} \quad (4.102)$$

We now have a set of five equations to solve for the five constants: from (4.87)

$$\lambda P_{00} = \mu_0 P_{10} + \mu_1 P_{01},$$

and (4.98), (4.100), (4.101), (4.102).

**Remark 4.17** In the derivation of the functional forms of  $f_{10}(x)$ ,  $f_{01}(x)$  the expressions

$$\mu_1 - \mu_0, \quad 2\mu_1 - \lambda - 2\mu_0, \quad \mu_1 - \lambda - \mu_0, \quad 2\mu_1 - \lambda + a$$

appear in various denominators. If any of these four expressions were equal to 0, the functional forms would have to be modified. The five equations used



to solve for the constants in the present model would have to be modified accordingly. In this monograph we emphasize the system-point level-crossing approach to derive model equations, and various techniques to solve them. However, there are many techniques to solve systems of integral equations, requiring additional study, outside the scope of the present volume. We give **numerical** solutions of the equations in several examples below.

**Remark 4.18** It would be interesting to explain the appearance of the immediately above expressions in the denominators. Does the system reduce to a particular queueing model when a denominator is equal to 0? For example, when  $\mu_1 - \mu_0 = 0$ , the  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$  system reduces to a standard M/M/2 model. In  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$  the only criterion necessary for stability is  $\lambda < 2\mu_1$ . What do these exceptional denominators mean with regard to physical models?

Another question is how to select a set of linearly independent equations to solve for the constants. Once a set of equations is derived, it can be checked for independence using matrix methods. But this amounts to trial and error. Is there a way to derive five independent equations directly? Taking derivatives may be the answer to this question.

**Example 4.7** We first give a mild numerical check on the five equations by letting  $\mu_1 - \mu_0 = 0$ . In this case  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$  reduces to a standard M/M/2 queue. We arbitrarily take

$$\lambda = 1, \mu_0 = 1.5, \mu_1 = 1.5.$$

Then  $a = -2.581139$ . The solution for the constants is

$$C_{10} = 0.0, P_{01} = .133333, P_{10} = .20, C_{01} = .333333, P_{00} = .50.$$

We compare this solution with that of the standard M/M/2 queue with  $\lambda = 1, \mu = 1.5$ . In M/M/2, the probability of an empty system is  $P_0 = 0.5$ . The probability of 1 customer in the system is indeed  $P_1 = 0.33333$ . The values match  $P_{00}$  and  $P_{10} + P_{01}$  in  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$  model, as expected.

Also, in  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$ , we see from (4.97) that

$$\begin{aligned} f_{01}(x) &= C_{01}e^{-(2\mu_1-\lambda)x} \\ &= \lambda P_1 e^{-(2\mu_1-\lambda)x} \\ &= 1 \cdot (0.33333)e^{-2x}, x > 0, \end{aligned}$$

since  $\mu_1 - \mu_0 = 0$  and  $C_{10} = 0$ .

**Example 4.8** Let  $\lambda = 1$ ,  $\mu_0 = 1.1$ ,  $\mu_1 = 2.21$ . These values preclude that any of the four above-mentioned denominators is 0. We get  $a = -2.715136$ . We solve the equations and obtain

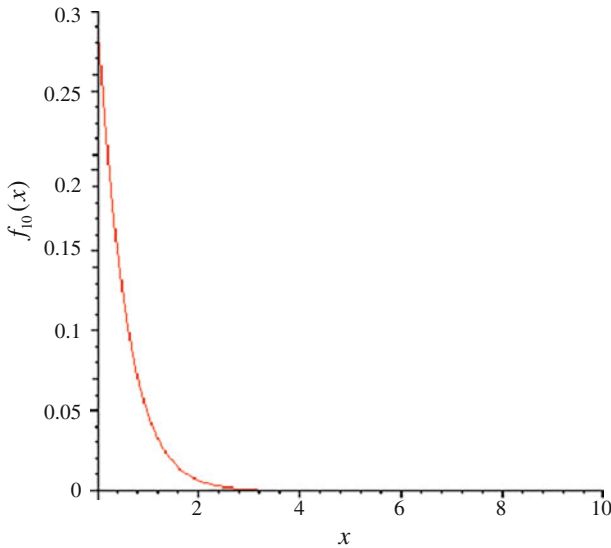
$$\begin{aligned} P_{00} &= .417715, & P_{10} &= 0.339103, & P_{01} &= 0.0202270, \\ C_{01} &= 0.022818, & C_{10} &= -0.322655. \end{aligned}$$

The functions  $f_{10}(x)$ ,  $x > 0$ , and  $f_{01}(x)$ ,  $x > 0$ , are linear combinations of exponentials,

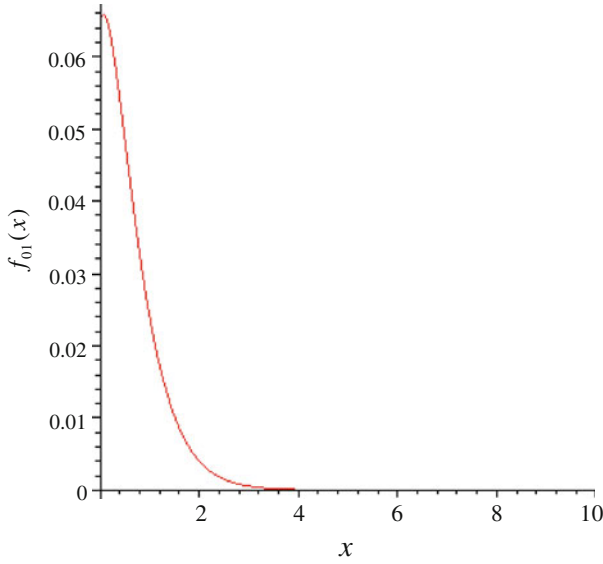
$$\begin{aligned} f_{10}(x) &= -0.322655e^{-2.715136x} + 0.617056e^{-2.2x}, \\ f_{01}(x) &= 0.505784e^{-2.2x} + 0.067831e^{-3.31x} \\ &\quad - 0.531504e^{-2.715136x} + 0.022818e^{-3.42x}. \end{aligned}$$

We substitute the values of  $P_{00}$ ,  $P_{10}$ ,  $P_{01}$ ,  $f_{10}(x)$ ,  $f_{01}(x)$  into the normalizer (4.91), and obtain 1; it checks.

The partial pdfs of wait  $f_{10}(x)$ ,  $f_{01}(x)$  and total pdf of wait  $f(x)$  are depicted in Figs. 4.7, 4.8, and 4.9 respectively.

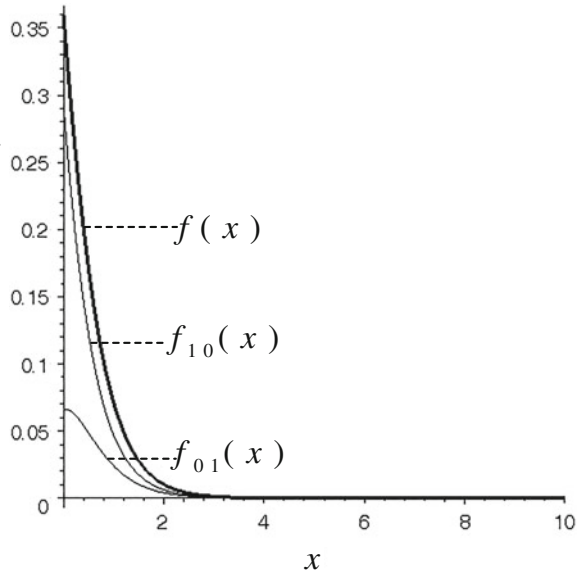


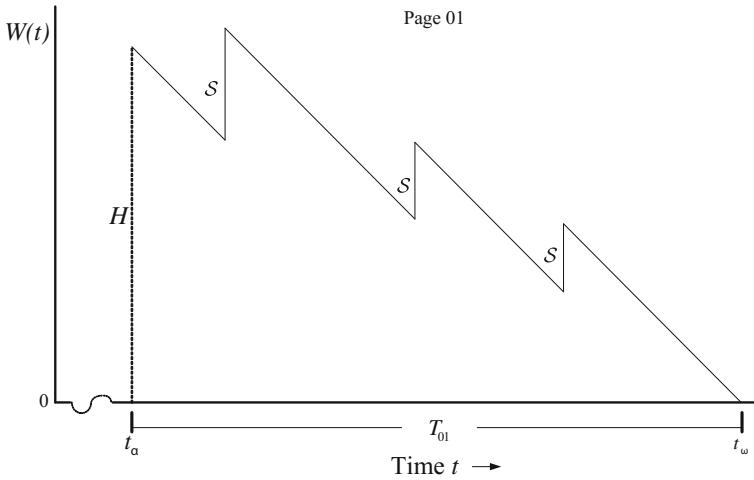
**Fig. 4.7** Partial pdf of wait  $f_{10}(x)$  in  $M/M/2/(\mu_0, \mu_1)$ ,  $(0, (0, \infty))$ .  $\lambda = 1$ ,  $\mu_0 = 1.1$ ,  $\mu_1 = 2.21$



**Fig. 4.8** Partial pdf of wait  $f_{01}(x)$  in  $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$ .  $\lambda = 1, \mu_0 = 1.1, \mu_1 = 2.21$

**Fig. 4.9** Total pdf of wait  
 $f(x) = f_{10}(x) + f_{01}(x)$   
 in  
 $M/M/2/(\mu_0, \mu_1), (0, (0, \infty))$ .  
 $\lambda = 1, \mu_0 = 1.1,$   
 $\mu_1 = 2.21$





**Fig. 4.10**  $T_{01} :=$  sojourn on page 01.  $t_\alpha :=$  start of  $T_{01}$ ,  $t_\omega :=$  end of  $T_{01}$ .  $S = \text{Exp}_{dis}^{2\mu_1}$ . (See Fig. 4.6.)

### 4.11.5 Expected Sojourn Time on a Page

Consider page 01. The SP can enter page 01 from discrete state  $(0, 01)$  or from page 10, due to a jump at an arrival (Fig. 4.6). It cannot enter directly from state  $(0, 10)$  at an arrival instant, since zero-wait arrivals are assigned rate  $\mu_0$  resulting in both servers being occupied with rate- $\mu_0$  customers; so any SP jump to a positive level must end on page 10.

In a sojourn on page 01, the first inter start-of-service depart time will be  $\stackrel{dis}{=} \text{Exp}_{\mu_0 + \mu_1}$ ; any other inter start-of-service depart times that follow while on page 01 will be  $\stackrel{dis}{=} \text{Exp}_{2\mu_1}$ . While the SP is on page 01, each departure will leave a rate- $\mu_1$  customer in the neighboring occupied server. Given that the SP enters page 01, its source state was  $(0, 01)$  with probability (using Bayes' rule)

$$q = \frac{P_{01}}{P_{01} + \int_{y=0}^{\infty} f_{10}(y)dy}.$$

Its source was composite state  $((0, \infty), 10)$  with probability

$$1 - q = \frac{\int_{y=0}^{\infty} f_{10}(y)dy}{P_{01} + \int_{y=0}^{\infty} f_{10}(y)dy}.$$

Let  $H$  denote the height above level 0 (ordinate) at which the SP enters page 01 (see Fig. 4.10). A sojourn on page 01 starts at level  $H$ , where

$$E(H|\text{source is } (0, 01)) = \frac{1}{\mu_0 + \mu_1},$$

and

$$E(H|\text{source is level } y \text{ on page } 10) = y + \frac{1}{\mu_0 + \mu_1}, \quad y > 0,$$

since the size of a jump from either source onto page 01 is  $\stackrel{\text{dis}}{=} \text{Exp}_{\mu_0 + \mu_1}$ .

Thus

$$E(H) = \frac{1}{\mu_0 + \mu_1} \cdot q + \left( \int_{y=0}^{\infty} \left( y + \frac{1}{\mu_0 + \mu_1} \right) f_{10}(y) dy \right) \cdot (1 - q).$$

From (4.93)  $f_{10}(y)$  is given by

$$f_{10}(y) = C_{10}e^{ay} + \lambda K_{10} P_{10} e^{-2\mu_0 y}, \quad y > 0,$$

and thus

$$\begin{aligned} E(H) &= \frac{1}{\mu_0 + \mu_1} \cdot q \\ &\quad + \left( \int_{y=0}^{\infty} \left( y + \frac{1}{\mu_0 + \mu_1} \right) (C_{10}e^{ay} + \lambda K_{10} P_{10} e^{-2\mu_0 y}) dy \right) \cdot (1 - q) \\ &= \frac{1}{\mu_0 + \mu_1} \cdot q + \left( \frac{1}{4} (4C_{10}\mu_0^2\mu_1 + 4C_{10}\mu_0^3 \right. \\ &\quad \left. + 3\lambda K_{10} P_{10} a^2 \mu_0 + \lambda K_{10} P_{10} a^2 \mu_1 \right. \\ &\quad \left. - 4C_{10}a\mu_0^2) / (a^2 \mu_0^2 (\mu_0 + \mu_1)) \right) \cdot (1 - q). \end{aligned} \quad (4.103)$$

Let  $T_{01}$  denote a sojourn time on page 01, i.e., the time from SP entrance until the first exit from page 01 thereafter. The only possible exit is due to a downcrossing of level 0 (Fig. 4.6). Thus

$$T_{01} = H + \sum_{i=1}^{N_H} \mathcal{B}_i$$

where  $N_H$  is the number of arrivals during time  $H$  and  $\mathcal{B}_i$  represents a busy period of an  $M/M/1$  queue with service rate  $2\mu_1$ , since both servers are busy with rate- $\mu_1$  customers. (See Sect. 3.4.12 and Fig. 3.6.) The expected busy period is obtained from (3.120) with  $2\mu_1$  substituted for  $\mu$ . Thus

$$E(\mathcal{B}_i) = \frac{1}{2\mu_1 - \lambda}, \quad i = 1, \dots, N_H.$$

The r.v.s  $N_H$  and  $\mathcal{B}_i$ ,  $i = 1, \dots, N_H$  are independent, since the  $\mathcal{B}_i$ s are i.i.d. each distributed as an  $M_\lambda/M_{2\mu_1}/1$  busy period. The expected sojourn time on page 01 is

$$\begin{aligned} E(T_{01}) &= E(H) + E\left(\sum_{i=1}^{N_H} \mathcal{B}_i\right) = E(H) + E(N_H)E(\mathcal{B}_i) \\ &= E(H) + \lambda E(H) \frac{1}{2\mu_1 - \lambda} = \frac{E(H)}{1 - \lambda/(2\mu_1)}, \end{aligned} \quad (4.104)$$

where  $E(H)$  is given in formula (4.103). It is noteworthy that  $T_{01}$  is distributed as the *busy period* in an  $M_\lambda/M_{2\mu_1}/1$  queue in which zero-wait arrivals obtain special service  $\underset{dis}{=} H$ , and positive-wait arrivals get service rate  $2\mu_1$ . This structure of  $T_{01}$  illustrates an interesting application, and the versatility, of the M/G/1 queue where zero-wait arrivals get exceptional service (see Sect. 3.6.1).

**Example 4.9** In Example 4.8 with  $\lambda = 1$ ,  $\mu_0 = 1.1$ ,  $\mu_1 = 2.21$ , we obtain

$$q = .111216, \quad 1 - q = .888784, \quad E(H) = 0.151416.$$

The expected sojourn time on page 01 is  $E(T_{01}) = 0.195689$ .

**Remark 4.19** Various questions arise regarding Example 4.9. What is the *proportion* of time that the SP spends circulating on page 01, page 10, or in the zero-wait states? Can this question be answered for a general M/M/c/ $(\mu_0, \mu_1)$ ,  $(0, (0, \infty))$  queue with  $c > 2$ ? If yes, then it would be straightforward to determine  $P_{00}$ . This would facilitate solving for all the zero-wait probabilities and the partial pdfs of wait.

## 4.12 M/M<sub>i</sub>/c with Reneging

Consider an M/M/c queue, with  $c \geq 2$  *distinguishable* servers having fixed exponential service rates  $\mu_i$ ,  $i = 1, \dots, c$ . Thus, the queue has *heterogeneous servers*. This model is denoted by M/M<sub>i</sub>/c. Using the notation for the generalized M/M/c model (Sects. 4.3, 4.4 and 4.5), let  $\{W(t), M(t)\}_{t \geq 0}$  denote the *system point process*, where  $W(t) :=$  *virtual wait* at time  $t$  and  $M(t) :=$  *system configuration* at time  $t$  (see Sect. 4.5). The set of possible exponential service rates is  $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_c\}$ . A new arrival receives one of those service rates, depending on which server it engages. We assume the

$\mu_i$ s are distinct. When some or all of the  $\mu_i$ s are equal, the analysis is similar with slight modification.

Assume zero-wait arrivals start service immediately (no balking). In general, the zero-wait server-assignment policy is arbitrary. When formulating equations for the zero-wait probabilities in a specific model, however, we must specify a zero-wait server-assignment policy (see Sect. 4.12.7 below).

### 4.12.1 Staying Function

Let  $\{\tau_n\}_{n=1,2,\dots}$  be the arrival times of customers  $C_n, n = 1, 2, \dots$ , respectively. Then  $W(\tau_n^-) \equiv W_n :=$  required wait before start of service of  $C_n$ .

Define

$$\theta_n = \begin{cases} 1 & \text{if } C_n \text{ stays for a full service} \\ 0 & \text{if } C_n \text{ reneges while waiting for service} \end{cases}, n = 1, 2, \dots$$

With respect to the steady-state statistical properties of the waiting time, this model is equivalent to one in which customers balk from joining the system at arrival instants, depending on their required wait before service, i.e., on their arrival-point  $W_n$ s. (See a sample path of  $\{W(t)\}_{t \geq 0}$  in Fig. 3.24 for a similar M/G/1 model with reneging.)

We define the *staying function*  $\bar{R}(\cdot)$  similarly as in Sect. 3.13.1. For each  $y \geq 0$ , define the *conditional probabilities*

$$\bar{R}(y) \equiv P(\theta_n = 1 | W_n = y), \quad R(y) \equiv P(\theta_n = 0 | W_n = y),$$

independent of  $n = 1, 2, \dots$ . Note that  $\bar{R}(0) = 1$ , and  $\bar{R}(y) + R(y) = 1, y \geq 0$ .

For each  $y \geq 0$ , given  $W_n = y, \theta_n$  has a Bernoulli distribution (e.g., p. 26 in [125]). The staying function  $\bar{R}(y)$  is the conditional probability of an arrival staying for a full service, given  $W_n = y$ . Its complement  $R(y)$  is the probability of an arrival reneging while in the waiting line, given  $W_n = y$ .

Using the foregoing definition,  $1 - \bar{R}(y), y \geq 0$ , is not necessarily a cdf.

We assume:  $\bar{R}(0) = 1; \bar{R}(y), y \geq 0$ , is monotone decreasing in the wide sense (i.e., not strictly monotone—it may be non-increasing);  $\bar{R}(y), y > 0$ , is bounded from below by 0.  $\bar{R}(y)$  may be continuous or piecewise continuous; it may be a step function.

Due to boundedness from below and monotonicity,  $\lim_{y \rightarrow \infty} \bar{R}(y)$  exists. Let

$$\lim_{y \rightarrow \infty} \bar{R}(y) = L, 0 \leq L \leq 1.$$

If  $\bar{R}(y) \equiv 1, y \geq 0$ , the model reverts to a standard M/M<sub>i</sub>/c queue with no renegeing; in that case  $L=1$  (see Sect. 3.13 and Theorem 3.8.)

### 4.12.2 System Configuration

The set of possible system configurations is

$$\mathbf{M} = \mathbf{M}_0 \cup \mathbf{M}_1 = \{\mathbf{m} | (m_1, m_2, \dots, m_c) | 0 \leq \sum_{i=1}^c m_i \leq c - 1\},$$

where  $m_i = \begin{cases} 1 & \text{if server } i \text{ is occupied} \\ 0 & \text{if server } i \text{ is idle} \end{cases}$ , just after a start of service in some server, since the configuration represents the service rates of those servers other than the one just occupied.

There are  $\binom{c}{j}$  configurations in which exactly  $j$  servers are occupied (i.e.,  $\sum_{i=1}^c m_i = j$ ). The total number of configurations in  $\mathbf{M}$  is

$$\sum_{j=0}^{c-1} \binom{c}{j} = 2^c - 1.$$

The number of configurations in  $\mathbf{M}_0 := \{\mathbf{m} | 0 \leq \sum_{i=1}^c m_i \leq c - 2\}$ , is  $2^c - 1 - c$ . The number of configurations in  $\mathbf{M}_1 := \{\mathbf{m} | \sum_{i=1}^c m_i = c - 1\}$  (border configurations), is  $c$ . (Recall that  $\mathbf{M}_1 = \mathbf{M}_b$ .)

### 4.12.3 State of System and Sample Path

#### State of System

Denote the state of the system as  $\{W(t), M(t)\}_{t \geq 0}$ , where  $W(t) \geq 0 :=$  virtual wait, and  $M(t) \in \mathbf{M} :=$  system configuration, at instant  $t$ .

#### Sample Path

Consider a sample path of  $\{(W(t), M(t))\}_{t \geq 0}$ . A sample-path diagram has  $2^c - 1$  lines corresponding to the zero-wait states  $(0, m), m \in \mathbf{M}$  (i.e.,  $W(t) = 0$ ); and  $c$  sheets (pages) corresponding to the positive-wait states  $(y, m), y > 0$  (i.e.,  $W(t) > 0$ ). (See Fig. 4.11 for the special case  $c = 2$ .)



Assume the system starts empty at  $t = 0$ . Initially, arriving customers wait 0, complete service and depart. Eventually customers in service accumulate until  $c - 1$  servers are occupied. Concurrently the SP moves among the  $2^c - 1 - c$  lines for the non-border zero-wait states. It resides on each such line for an exponentially distributed time, making transitions from line to line. Various states unfold until the SP ends up on one of the  $c$  border lines.

All zero-wait arrivals stay for full service (no balking). Assume that a new arrival  $C_\tau$  finds  $c - 1$  servers occupied (SP on a border line). Then  $C_\tau$  waits 0, and starts service in the single idle server. At  $\tau^-$  the configuration is some  $\mathbf{m} \in \mathbf{M}_b$ . At instant  $\tau$  all  $c$  servers are occupied. The SP jumps at instant  $\tau$  to one of the  $c$  sheets, depending on which service will finish first. The probability that server  $k$  will finish first is  $\mu_k/\mu$  where  $\mu := \mu_1 + \dots + \mu_c$ . The SP will be at a height  $\stackrel{dis}{=} \text{Exp}_\mu$ , since the inter start-of-service depart time  $\mathcal{S}$  is the *minimum* of  $c$  independent exponentially distributed r.v.s with rates  $\mu_1, \dots, \mu_c$ , due to the memoryless property.

Let  $\mathbf{m}_{\bar{i}}$  denote a *border* configuration such that the rate- $\mu_i$  server (i.e., server  $i$ ) is *idle* (see Remark 4.20). In configuration  $\mathbf{m}_{\bar{i}}$ ,  $m_j = 1$ , if  $j \neq i$ , and  $m_i = 0$ , i.e.,

$$m_1 + \dots + m_{i-1} + 0 + m_{i+1} + \dots + m_c = c - 1.$$

At time  $\tau$  the SP will end up at a positive height on page  $\mathbf{m}_{\bar{k}}$  with probability  $\mu_k/\mu$ ,  $k = 1, \dots, c$ .

**Remark 4.20** We use the notation  $\bar{i}$  to shorten the representation of  $\mathbf{m}$  if  $c$  is large. If  $c$  is small, e.g.  $c = 3$ , we can use notation like 100, 010, 001, 110, 101, 011. If  $c = 2$ , we can use 01, 10—see Sect. 4.12.8.

#### 4.12.4 Zero-Wait Probabilities

Let  $P_n$ ,  $n = 0, \dots, c - 1$  denote the steady-state probability of  $n$  customers in the system at an arbitrary point in time. Let  $P_{n,m}$  denote the probability that there are  $n$  customers in the system and the configuration is  $\mathbf{m} \in \mathbf{M}$ . There are  $\binom{c}{n}$  configurations such that  $\sum_{i=1}^c m_i = n$ . Let

$$\mathbf{M}_n = \{\mathbf{m} \mid \sum_{i=1}^c m_i = n\}.$$

Thus

$$P_n = \sum_{m \in M_n} P_{n,m}, n = 0, \dots, c - 1.$$

Due to Poisson arrivals  $P_n$  is the probability that an arrival waits 0 and “sees”  $n$  other customers in service just before it starts service (using PASTA, e.g., [145]).

**Remark 4.21** For the **zero-wait** states, a configuration specifies the service rates in the servers at an arbitrary time point. Due to Poisson arrivals, this is the same as the service rates **just before** an arrival. It is also the same as the service rates in the **other** servers **just after** an arrival starts service in an available server.

The probability of a zero wait is denoted by  $F(0)$ , where

$$F(0) = \sum_{n=0}^{c-1} P_n = \sum_{n=0}^{c-1} \sum_{m \in M_n} P_{n,m}. \quad (4.105)$$

#### 4.12.5 Positive-Wait PDF and CDF

For the **positive-wait** states, a configuration defines the service rates in the **other** servers **just after** start of service.

Let  $f_m(x)$ ,  $x > 0$ , denote the *partial* pdf of wait for page  $m \in M_b$ . Denote the marginal pdf for the *cover* as

$$f(x) = \sum_{m \in M_1} f_m(x), x > 0.$$

The *total* pdf of wait is  $\{P_0, f(x), x > 0\}$ . The cdf of wait is  $F(x) = F(0) + \int_{y=0}^x f(y)dy$ ,  $x \geq 0$ , where  $F(0)$  is defined in (4.105). The normalizing condition is

$$\lim_{x \rightarrow \infty} F(x) = F(0) + \int_{y=0}^{\infty} f(y)dy = 1. \quad (4.106)$$

#### 4.12.6 Equations for Positive-Wait PDFs

A key assumption of this model is that each positive-wait arrival reneges from the waiting line with probability  $R(y)$ , and stays for complete service

with probability  $\bar{R}(y) (=1 - R(y))$ , where  $y \geq 0$  is the required wait before service.

**Equation for Total PDF  $f(x)$**  We first derive an integral equation for  $f(x)$ , the total pdf of wait of stayers (who wait and receive a full service), namely,

$$f(x) = \lambda P_{c-1} e^{-\mu x} + \lambda \int_{y=0}^x e^{-\mu(x-y)} \bar{R}(y) f(y) dy, \quad x > 0, \quad (4.107)$$

directly using the sample path, as follows (see Fig. 4.11).

**Explanation of Equation (4.107)** On the left side,  $f(x)$  is the *total SP downcrossing rate* of level  $x$  over all  $c$  sheets, *projected onto the “cover”*. On the right side, since all zero-wait arrivals stay for full service ( $\bar{R}(0) = 1$ ), the term  $\lambda P_{c-1} e^{-\mu x}$  is the total SP upcrossing rate of level  $x$  due to jumps starting at level 0 (i.e., line 0) of any of the  $c$  sheets (from border states  $\{(0, \mathbf{m}_{\bar{i}})\}$ ,  $i = 1, \dots, c$ ), at arrival instants. These jumps have size  $\mathcal{S} = \text{Exp}_{dis} \mu$  ( $= \min_{i=1, \dots, c} \{\text{Exp}_{\mu_i}\}$ ). The term  $\lambda \int_{y=0}^x e^{-\mu(x-y)} \bar{R}(y) f(y) dy$  is the rate at which the SP upcrosses level  $x$  due to jumps starting at levels  $y \in (0, x)$ , on any page, at arrival instants of stayers. The right side is, therefore, the total SP upcrossing rate of level  $x$ . Rate balance across  $x$  yields (4.107).

Comparing (4.107) with Eq. (3.211) implies that the solution of (4.107) is

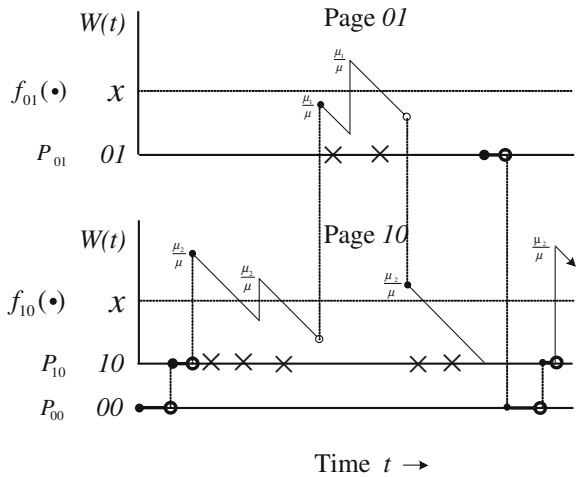
$$f(x) = \lambda P_{c-1} e^{-\left(\mu x - \lambda \int_{y=0}^x \bar{R}(y) dy\right)}, \quad x > 0, \quad (4.108)$$

where  $\mu = \sum_{i=1}^c \mu_i$  and  $P_{c-1} = \sum_{i=1}^c P_{c-1, \mathbf{m}_{\bar{i}}}$ .

**Equations for Partial PDFs  $f_{\bar{i}}(x)$ ,  $x > 0$ ,  $i = 1, \dots, c$**  We now obtain integral equations for the pdfs  $f_{\bar{i}}(x)$ ,  $x > 0$ , on the  $c$  sheets (see Fig. 4.11); they are

$$\begin{aligned} f_{\bar{i}}(x) + \lambda \left(1 - \frac{\mu_i}{\mu}\right) \int_{y=x}^{\infty} \bar{R}(y) f_{\bar{i}}(y) dy \\ = \lambda \frac{\mu_i}{\mu} P_{c-1} e^{-\mu x} + \lambda \frac{\mu_i}{\mu} \int_{y=0}^x e^{-\mu(x-y)} \bar{R}(y) f(y) dy \\ + \lambda \frac{\mu_i}{\mu} \int_{y=x}^{\infty} \bar{R}(y) (f(y) - f_{\bar{i}}(y)) dy, \quad i = 1, \dots, c. \end{aligned} \quad (4.109)$$

**Fig. 4.11** Sample path of  $\{W(t), M(t)\}_{t \geq 0}$  in  $M/M_i/2$ , where zero-wait arrivals join, and positive-wait arrivals may renege from the waiting line. Times marked  $\times$  indicate arrivals that renege (do not contribute to the limiting pdf of wait)



**Explanation of Equation (4.109)** On the left side,  $f_{\bar{i}}(x)$  is the SP exit rate from composite state  $((x, \infty), \bar{i})$  due to SP downcrossings of level  $x$ ; term  $\lambda(1 - \frac{\mu_i}{\mu}) \int_{y=x}^{\infty} \bar{R}(y) f_{\bar{i}}(y) dy$  is the SP rate of jumps out of  $((x, \infty), \bar{i})$  into the composite states  $((x, \infty), \bar{j})$ ,  $j \neq i$ , on other sheets. On the right side, the first two terms are SP entrance rates into  $((x, \infty), \bar{i})$  due to jumps starting at level-0 border states, and jumps starting at levels  $y \in (0, x)$  on any sheet, respectively (recall  $f(y) = \sum_{i=1}^c f_{\bar{i}}(y)$ ). The third term is the SP entrance rate into  $((x, \infty), \bar{i})$  due to jumps starting in  $\cup_{j \neq i} ((x, \infty), \bar{j})$ . Rate balance of SP exits and entrances of  $((x, \infty), \bar{i})$  yields (4.109).

**Solution of Equation (4.109)** We obtain the solution of (4.109) in terms of the solution for  $f(x)$ , which is given in formula (4.108), using the following Proposition.

**Proposition 4.1** *The partial pdf is given by*

$$f_{\bar{i}}(x) = \frac{\mu_i}{\mu} f(x), x > 0, i = 1, \dots, c. \tag{4.110}$$

**Proof** Substitute  $\frac{\mu_i}{\mu} f(x)$  for  $f_{\bar{i}}(x)$  in Eq. (4.109), and cancel like terms. The proposition is true if and only if the following is an identity:

$$\begin{aligned}
& \frac{\mu_i}{\mu} f(x) + \lambda \int_{y=x}^{\infty} \frac{\mu_i}{\mu} \bar{R}(y) f(y) dy \\
&= \lambda \frac{\mu_i}{\mu} P_{c-1} e^{-\mu x} + \lambda \frac{\mu_i}{\mu} \int_{y=0}^x e^{-\mu(x-y)} \bar{R}(y) f(y) dy \\
&\quad + \lambda \frac{\mu_i}{\mu} \int_{y=x}^{\infty} \bar{R}(y) f(y) dy, \quad x > 0, \tag{4.111}
\end{aligned}$$

if and only if

$$f(x) = \lambda P_{c-1} e^{-\mu x} + \lambda \int_{y=0}^x e^{-\mu(x-y)} \bar{R}(y) f(y) dy, \quad x > 0, \tag{4.112}$$

is an identity. Equation (4.112) is identical to Eq. (4.107). Hence the Proposition is true. ■

### Exponential Staying Function

Consider an exponential staying function,  $\bar{R}(x) := e^{-rx}$ ,  $r > 0$ ,  $x \geq 0$ . (Note that  $0 < e^{-rx} \leq 1$ , and is strictly decreasing on  $(0, \infty)$ , satisfying the definition of staying function.) The total pdf  $f(x)$  is now obtained by substituting  $e^{-ry}$  for  $\bar{R}(y)$  in (4.108), which substituted into (4.110), gives

$$f_{\bar{i}}(x) = \lambda \frac{\mu_i}{\mu} e^{\frac{\lambda}{r}} P_{c-1} e^{-\mu x - \frac{\lambda}{r} e^{-rx}}, \quad x > 0, \quad i = 1, \dots, c. \tag{4.113}$$

We shall solve an M/M<sub>i</sub>/2 model using  $\bar{R}(x) := e^{-rx}$ ,  $r > 0$ ,  $x \geq 0$ , in Sect. 4.12.8 below.

#### 4.12.7 Equations for Zero-Wait Probabilities

Assume that the zero-wait server assignment policy is: arrivals that find  $k$  available servers,  $1 \leq k \leq c$ , get served by a particular available server with probability  $1/k$ . (Other policies are also viable, e.g., the arrival gets served by the lowest-numbered available server, or by the fastest-available service rate, etc.) Using the principle *SP exit rate = SP entrance rate* for the zero-wait states, we obtain the equations (notation explained below)

$$\begin{aligned}
 (\lambda + \mu - \mu_i)P_{c-1,\bar{i}} &= f_{\bar{i}}(0) + \frac{\lambda}{2} \sum_{j \in J_i} P_{c-2,\bar{i}\bar{j}}, \quad i = 1, \dots, c, \\
 (\lambda + \mu - \mu_i - \mu_j)P_{c-2,\bar{i}\bar{j}} &= \mu_j P_{c-1,\bar{i}} + \mu_i P_{c-1,\bar{j}} \\
 &\quad + \frac{\lambda}{3} \sum_{k \in J_{ij}} P_{c-3,\bar{i}\bar{j}\bar{k}}, \quad j = 1, \dots, c, \\
 &\quad \dots \\
 (\lambda + \mu_i)P_{1,i} &= \sum_{k \neq i=1}^c \mu_k P_{2,ik} + \frac{\lambda}{c} P_{00}, \quad i = 1, \dots, c, \\
 \lambda P_{00} &= \sum_{i=1}^c \mu_i P_{1,i}.
 \end{aligned} \tag{4.114}$$

**Notation in equations (4.114)** In the first  $c$  equations, the index  $j$  of the sum takes values in  $J_i := \{j | j = 0, \dots, c, j \neq i\}$ , and the subscript  $\bar{i}\bar{j}$  means both servers  $i$  and  $j$  are idle. In the second set of  $\binom{c}{2}$  equations, the index  $k$  of the sum takes values in  $J_{ij} = \{k | k = 0, \dots, c, k \neq i, k \neq j\}$ , and the subscript  $\bar{i}\bar{j}\bar{k}$  means all three servers  $i, j$  and  $k$  are idle. The row of dots “...” indicates similar rate balance equations for  $P_{c-3,\bar{i}\bar{j}\bar{k}}, \dots, P_{2,}$ . In the second last equation, for  $P_{1,i}$ , on the right side  $P_{2,ik}$  denotes the probability of two units in the system, in servers  $i$  and  $k$  having service rates  $\mu_i$  and  $\mu_k$  respectively.

We solve Eq. (4.114) explicitly in Sect. 4.12.8 below for M/M<sub>i</sub>/2, in order to convey some characteristics of the solution.

### 4.12.8 Solution for M/M<sub>i</sub>/2 with Reneging

**Notation** When there is a small number of servers we can use an alternative, perhaps more familiar notation. If  $c = 2$ , there are two sheets corresponding to configurations  $\bar{1}$  and  $\bar{2}$ , which we now replace by 01 and 10 respectively. Thus, configuration 01 means server 1 is available and server 2 is occupied; configuration 10 means server 2 is available and server 1 is occupied.

Applying formula (4.113), the partial pdfs of wait are now denoted by

$$\begin{aligned}
 f_{10}(x) &= \lambda \frac{\mu_2}{\mu} e^{\frac{\lambda}{r}} P_1 e^{-\mu x - \frac{\lambda}{r} e^{-rx}}, \quad x > 0, \\
 f_{01}(x) &= \lambda \frac{\mu_1}{\mu} e^{\frac{\lambda}{r}} P_1 e^{-\mu x - \frac{\lambda}{r} e^{-rx}}, \quad x > 0.
 \end{aligned} \tag{4.115}$$

The marginal (“total”) pdf of wait is

$$f(x) = f_{10}(x) + f_{01}(x) = \lambda e^{\frac{\lambda}{r}} P_1 e^{-\mu x - \frac{\lambda}{r} e^{-rx}}, \quad x > 0. \tag{4.116}$$

The zero-wait probabilities are  $P_{1,\bar{i}}$ ,  $i = 1, 2$ , and  $P_{00}$ ; using the alternative notation we have

$$\begin{aligned} P_1 &= P_{1,\bar{2}} + P_{1,\bar{1}} = P_{10} + P_{01}, \\ P_0 &= P_{00} + P_{1,\bar{2}} + P_{1,\bar{1}} = P_0 + P_{10} + P_{01} \\ &= P_{00} + P_1. \end{aligned}$$

The rate-balance equations for the zero-wait probabilities are

$$\begin{aligned} (\lambda + \mu_1)P_{10} &= \frac{\lambda}{2}P_{00} + f_{10}(0), \\ (\lambda + \mu_2)P_{01} &= \frac{\lambda}{2}P_{00} + f_{01}(0), \\ \lambda P_{00} &= \mu_1 P_{10} + \mu_2 P_{01}. \end{aligned} \tag{4.117}$$

Substituting for  $f_{10}(0)$ ,  $f_{01}(0)$  from (4.115), we rewrite the equations in (4.117) as

$$\begin{aligned} (\lambda + \mu_1)P_{10} &= \frac{\lambda}{2}P_{00} + \lambda \frac{\mu_2}{\mu} P_1, \\ (\lambda + \mu_2)P_{01} &= \frac{\lambda}{2}P_{00} + \lambda \frac{\mu_1}{\mu} P_1, \\ \lambda P_{00} &= \mu_1 P_{10} + \mu_2 P_{01}. \end{aligned} \tag{4.118}$$

The solution of (4.118) in terms of  $P_{00}$  is

$$\begin{aligned} P_{01} &= \frac{\lambda}{2\mu_2} P_{00}, \\ P_{10} &= \frac{\lambda}{2\mu_1} P_{00}, \\ P_1 &= \frac{\lambda(\mu_1 + \mu_2)}{2\mu_1\mu_2} P_{00} = \frac{\lambda\mu}{2\mu_1\mu_2} P_{00}. \end{aligned} \tag{4.119}$$

The normalizing condition

$$P_{00} + P_1 + \int_{x=0}^{\infty} f(x)dx = 1,$$

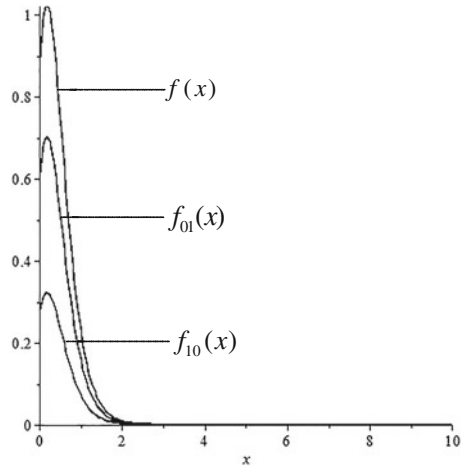
yields

$$P_{00} = \left( 1 + \frac{\lambda(\mu_1 + \mu_2)}{2\mu_1\mu_2} + \frac{\lambda(\mu_1 + \mu_2)}{2\mu_1\mu_2} \lambda e^{\frac{\lambda}{r}} \int_{x=0}^{\infty} e^{-\mu x - \frac{\lambda}{r} e^{-rx}} dx \right)^{-1}. \tag{4.120}$$

The analytic solution comprises the results in (4.120), (4.119), (4.116) and (4.115).

**Example 4.10** We present a numerical example for the M/M<sub>i</sub>/2 queue with reneging allowed from the waiting line (see Fig. 4.12). Let

**Fig. 4.12** Plot of  $f(x)$ ,  $f_{01}(x)$  ( $= f_{\bar{1}}(x)$ ),  $f_{10}(x)$  ( $= f_{\bar{2}}(x)$ ), in Example 4.10



$$\lambda = 5.2, \mu_1 = 2.4, \mu_2 = 1.1, \mu = \mu_1 + \mu_2 = 3.5, r = 2.1.$$

Then

$$\int_{x=0}^{\infty} e^{-\mu x - \frac{\lambda}{r} e^{-rx}} dx = 0.074741,$$

and

$$P_{00} = 0.049059, \quad P_{01} = 0.115958, \quad P_{10} = 0.053147,$$

$$P_1 = 0.169105,$$

$$F(0) = P_{00} + P_1 = 0.218164,$$

$$F(\infty) = F(0) + \lambda e^{\frac{\lambda}{r}} P_1 \int_{x=0}^{\infty} e^{-\mu x - \frac{\lambda}{r} e^{-rx}} dx$$

$$= 0.218164 + 0.781836 = 1.0,$$

$$f(x) = \lambda e^{\frac{\lambda}{r}} P_1 e^{-\mu x - \frac{\lambda}{r} e^{-rx}} = 10.461 \cdot e^{-3.5x - 2.476e^{-2.1x}},$$

$$f_{01}(x) = \frac{\mu_1}{\mu} f(x) = 7.173 \cdot e^{-3.5x - 2.476e^{-2.1x}},$$

$$f_{10}(x) = \frac{\mu_2}{\mu} f(x) = 3.288 \cdot e^{-3.5x - 2.476e^{-2.1x}}.$$

**Remark 4.22** In  $M/M_i/c$  with renegeing from the waiting line allowed, we can generalize the staying function  $\bar{R}(x)$ ,  $x \geq 0$ . For example,  $\bar{R}(x)$  may depend on the server that would be occupied by an arrival, i.e., on the system configuration at the arrival instant. We may then use the notation  $\bar{R}_{\bar{i}}(x)$ . Thus  $\bar{R}_{\bar{i}}(x)$  may depend on, not only customer required wait before service, but also



on customer attraction or aversion to the “target” server. A natural question arises. Can this model be modified to study attraction or aversion in natural processes such as: electrically charged particles approaching an electrically charged or magnetized environment; asteroids approaching a planet; particles adhering or falling away from a surface; laser pulses affecting cells containing certain chemicals in biological or medical applications; etc.?

### 4.12.9 Stability Condition

Consider the M<sub>λ</sub>/M<sub>i</sub>/c ( $c \geq 2$ ) queue with heterogeneous servers having rates  $\mu_1, \dots, \mu_c$  in which reneging depending on required wait is allowed before service begins. Let the staying function  $\bar{R}(x), x \geq 0$ , be monotone decreasing (includes non-increasing), let  $\bar{R}(0) = 1$  (no balking upon arrival), and assume  $0 \leq \bar{R}(x) \leq 1, x \geq 0$ . Let  $L = \lim_{x \rightarrow \infty} \bar{R}(x)$ , which exists by monotonicity and boundedness. The ideas in Theorem 3.8 also apply in the M/M<sub>i</sub>/c environment, as follows.

**Theorem 4.10** In M<sub>λ</sub>/M<sub>i</sub>/c ( $c \geq 2$ ) with reneging from the waiting line allowed, as described immediately above, a necessary and sufficient condition for stability is

$$\begin{aligned} \lambda &< \frac{\mu}{L} \text{ if } 0 < L \leq 1, \\ \lambda &< \infty \text{ if } L = 0, \end{aligned}$$

where  $\mu = \sum_{i=1}^c \mu_i$ .

**Proof** The proof is similar to that of Theorem 3.8. The alternative proof given there, Remark 3.31 and Fig. 3.28 also apply for the present M/M<sub>i</sub>/c queue with reneging, upon substituting  $\mu = \sum_{i=1}^c \mu_i$ . ■

### 4.13 Discussion

We can use LC to analyze a vast array of additional M/M/c models. We mention only a few examples.

LC has been applied to M/M/c queues in which customers receive simultaneous service from a random number of servers. The original source for such queueing models is the Ph.D. thesis of L. Green [81, 82]. An LC analysis, motivated by the work of L. Green, is given in [38].

LC has been applied to M/M/c with bounded system time (wait + service). An arrival balks upon arrival if its system time would exceed an upper bound  $K$ , e.g., a system manager informs an arrival of the current expected system time (see Example 1, p. 44 in [52]). This generalizes variant 2 of the M/G/1 model discussed above in Sect. 3.6. It is straightforward to apply LC to analyze an M/M/c model analogous to variant 1 in Sect. 3.16. In that model customers renege from *service* if their *age* in the system (*elapsed system time*) reaches  $K$ . Similar remarks apply to M/M/c where the actual waits are bounded by  $K$  (as in variant 3 in Sect. 3.16). In that case the workload can exceed  $K$ . We can develop an expression for the tail of the steady-state pdf of workload, from its integral equation.

LC can be used to analyze a variety of M/M/c queues with server vacations; priorities; and many others.