

# Warteschlangentheorie und Callcenter

Vortrag im Rahmen der Lehrerfortbildung  
„Stochastik und Matrizen: von Markov-Ketten bis zu Callcentern“  
23. September 2009

Dr. Alexander Herzog,  
Institut für Mathematik, TU Clausthal

# 1 Warteschlangenmodelle

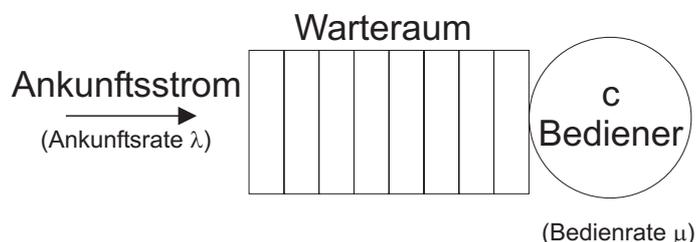


Abbildung 1: Allgemeines Warteschlangenmodell

Das Grundmodell in Abbildung 1 besteht aus folgenden Komponenten:

- **Ankunftsstrom:**

Der Ankunftsstrom wird durch die zufälligen Ankunftszeitpunkte der Kunden am System charakterisiert und wird durch die Verteilung der Zeiten zwischen zwei Kundenankünften, den sogenannten **Zwischenankunftszeiten**, beschrieben. Es sei  $I$  die Zufallsvariable der Zwischenankunftszeiten, dann wird  $\lambda := \frac{1}{\mathbf{E}[I]}$  die **Ankunftsrate** pro Zeiteinheit genannt. Es wird angenommen, dass die aufeinanderfolgenden Zwischenankunftszeiten stochastisch unabhängig und identisch verteilt sind.

- **Wartezimmer:**

Vor dem Bedienschalter befindet sich ein Wartezimmer, in dem sich die eingetroffenen Kunden aufhalten, bevor sie bedient werden. Der Wartezimmer kann eine endliche oder eine unendliche Kapazität besitzen. Besitzt der Wartezimmer nur eine endliche Kapazität, so werden eintreffende Kunden abgewiesen, sobald alle Warteplätze belegt sind.

- **Bedienprozess:**

An dem Bedienschalter arbeiten eine gewisse Anzahl an Bedienern parallel nebeneinander. Die **Bedienzeiten** der aufeinanderfolgenden Kunden bilden eine Folge stochastisch unabhängiger und identisch verteilter Zufallsvariablen mit Verteilungsfunktion  $F_S$  ( $S$  von englisch *Service*) und Erwartungswert  $\mathbf{E}[S]$ .  $\mu := \frac{1}{\mathbf{E}[S]}$  entspricht der **Bedienrate** pro Zeiteinheit eines Bedieners.

- **Bedienregeln:**

Über die sogenannten Bedienregeln wird festgelegt, nach welchem Prinzip die Kunden aus dem Wartezimmer zu den Bedienern geleitet werden. Da in den meisten Fällen der Ankunftszeitpunkt für die Priorität des Kunden im Wartezimmer verantwortlich ist, spricht man auch davon, dass sich vor dem Bedienschalter eine Warteschlange bildet.

## 1 Definition (Auslastung, Arbeitslast):

In einem System mit  $c \in \mathbb{N}$  Bedienern definiert  $\rho := \frac{\lambda}{c\mu}$  die Systemauslastung, wobei  $\mu$  die Bedienrate eines einzelnen Bedieners pro Zeiteinheit darstellt. Außerdem wird  $a := \frac{\lambda}{\mu}$  die **Arbeitslast** genannt.

Die Arbeitslast wird auch häufig zu Ehren von AGNER KRARUP ERLANG, dem Begründer der Warteschlangentheorie, mit der Einheit **Er1** (für „Erlang“) versehen, obwohl sie als Quotient zweier Raten einheitenlos ist.

## 2 Bezeichnung (Zustandswahrscheinlichkeit):

Im Folgenden sei  $\pi_n$ ,  $n \in \mathbb{N}$ , die Wahrscheinlichkeit dafür, dass sich in einem Warteschlangensystem im stationären Zustand  $n$  Kunden im System, d.h. in der Warteschlange oder in Bedienung, befinden.

### 3 Definition (mittlere Anzahl an Kunden im System):

Die mittlere Anzahl an Kunden im System in einem Warteschlangensystem wird als

$$\mathbf{E}[N] := \sum_{n=0}^{\infty} n\pi_n$$

definiert. Besteht ein Warteschlangensystem aus  $c$  Bedienern, so definiert man die mittlere Anzahl an Kunden in der Warteschlange als

$$\mathbf{E}[N_Q] := \sum_{n=c+1}^{\infty} (n - c)\pi_n.$$

### 4 Bezeichnung (mittlere Wartezeit, mittlere Verweildauer):

In einem Warteschlangensystem wird mit  $W$  die Zufallsvariable der Wartezeit der Kunden beschrieben und mit  $V$  die Zufallsvariable der Verweilzeit der Kunden im System, d.h. der Summe aus Warte- und Bedienzeit.

### 5 Satz (JOHN D. C. LITTLE, 1961):

Für Warteschlangensysteme mit allgemeinen Zwischenankunfts- und Bedienzeiten sowie  $c$  parallelen Bedienern gilt im stationären Zustand folgender Zusammenhang zwischen der mittleren Warteschlangenlänge  $\mathbf{E}[N_Q]$  und der mittleren Wartezeit der Kunden  $\mathbf{E}[W]$ :

$$\mathbf{E}[N_Q] = \lambda \mathbf{E}[W]$$

und folgender Zusammenhang zwischen der mittleren Anzahl an Kunden im System  $\mathbf{E}[N]$  und der mittleren Verweilzeit der Kunden  $\mathbf{E}[V]$ :

$$\mathbf{E}[N] = \lambda \mathbf{E}[V].$$

## 2 Die Kendall-Notation

Viele Warteschlangenmodelle lassen sich nach der sogenannten Kendall-Notation (1953 von DAVID GEORGE KENDALL eingeführt) klassifizieren. Nach der Kendall-Notation wird einem Warteschlangensystem eine Buchstabenkombination der folgenden Form zugeordnet:

**A / S / c / [C] [+ W] / [p] / [D]**

Eckige Klammern deuten dabei an, dass der jeweilige Parameter optional ist. Die Platzhalter „A“, „S“, „c“, „C“, „W“, „p“ und „D“ haben dabei diese Bedeutungen:

- „A“ gibt die Verteilung der Zwischenankunftszeiten (arrival) an.
- „S“ gibt die Verteilung der Bedienzeiten (service) an.
- „c“ gibt die Anzahl der identischen Bediener an. Es ist  $c \in \mathbb{N}$ . Übliche Werte sind „1“ oder „c“ um anzudeuten, dass eine feste Anzahl  $c \geq 1$  an Bedienern zur Verfügung steht.
- „C“ gibt die maximale Anzahl an Kunden im System an (wartende Kunden sowie in Bedienung befindliche Kunden). Fehlt diese Angabe, so wird  $C = \infty$  angenommen. Übliche Werte sind „c“, „N“ (mit  $N > c$ ) und „∞“.
- **W** gibt die Wartezeittoleranz-Verteilung der Kunden an. Fehlt diese Angabe, so besitzen alle Kunden eine unendliche Wartezeittoleranz und brechen einen Wartevorgang nie ab.

- „**P**“ gibt die Populationsgröße, d.h. die Gesamtzahl an Kunden, die am System ankommen können, an. Fehlt diese Angabe, so wird  $p = \infty$  angenommen.
- „**D**“ beschreibt die Bedienregel (*service discipline*). Üblich sind hier insbesondere folgende Werte:
  - **FIFO** (*First in first out*) oder auch **FCFS** (*First come first served*): Bedienung der Kunden in der Ankunftsreihenfolge.
  - **LIFO** (*Last in first out*) oder auch **LCFS** (*Last come first served*): Bedienung der Kunden in umgekehrter Ankunftsreihenfolge, d.h. es wird immer der zuletzt eingetroffene Kunde als nächstes bedient.
  - **SPT** (*Shortest processing time*): Es wird immer derjenige wartende Kunde als nächstes bedient, dessen Bedienzeit am kürzesten ist.
  - **RANDOM** (zufällig) oder auch **SIRO** (*Serve in random order*): Der nächste zu bedienende Kunde wird jeweils zufällig ausgewählt.

Wird keine Bedienregel angegeben, so wird stets FIFO angenommen.

Für die unterschiedlichen Verteilungstypen existieren folgende allgemein übliche Kürzel:

Symbol	Verteilung
$M$	Exponentialverteilung ( <u>M</u> arkov-Eigenschaft)
$D$	konstante Zwischenankunfts- bzw. Bedienzeiten ( <u>d</u> eterministisch)
$E_k$	Erlang- $k$ -Verteilung
$G$	unbekannte Verteilung ( <u>g</u> enerell), es sind nur Kenngrößen bekannt
$H_k$	Hyperexponentialverteilung (Linearkombination von $k$ Exponentialverteilungen)
$PH$	Phasentypverteilung (Verteilung der Zeit, in der ein Markov-Prozess einen absorbierenden Zustand erreicht)

### 3 Unabhängigkeit der Kundenankünfte

Wenn sich die Kunden in Bezug auf Ihre Ladenbesuche nicht absprechen und auch nicht z.B. busladungsweise eintreffen, kann davon ausgegangen werden, dass die Kunden unabhängig von einander eintreffen. Betrachtet man ein solches System zum einem Zeitpunkt  $t_0 \in \mathbb{R}_0^+$ , so ist die Wahrscheinlichkeit, dass innerhalb der nächsten 5 Minuten ein Kunde eintrifft, unabhängig davon, ob innerhalb der letzten 5 Minuten ein Kunde eingetroffen ist oder nicht. Als bedingte Wahrscheinlichkeit formuliert bedeutet dies:

$$P(X \leq t_1 | X > t_0) = P(X \leq t_1 - t_0), \quad (1)$$

wobei  $t_1 > t_0$  sei. Im einzelnen ist  $P(X \leq t_1 | X > t_0)$  die Wahrscheinlichkeit, dass die Ankunft eines Kunden vor oder zu dem Zeitpunkt  $t_1$  stattfindet ( $X \leq t_1$ ), unter dem Vorauswissen, dass bis zum Zeitpunkt  $t_0$  noch kein Kunde eingetroffen ist ( $X > t_0$ ) und  $P(X \leq t_1 - t_0)$  ist die Wahrscheinlichkeit, dass ein Kunde in der Zeitdifferenz zwischen  $t_0$  und  $t_1$  eingetroffen ist (ohne ein Vorauswissen über die Zeit zwischen 0 und  $t_0$ ).

#### 6 Definition (Bedingte Wahrscheinlichkeiten):

Es sei  $P(B) > 0$ , d.h. das Ereignis  $B$  trete mit einer Wahrscheinlichkeit größer als 0 ein. Dann sei  $P(A|B)$  die bedingte Wahrscheinlichkeit dafür, dass  $A$  eintritt, wenn bekannt ist, dass  $B$  eingetreten ist. Es gilt

$$P(A|B) := \frac{P(A \cap B)}{P(B)},$$

wobei  $P(A \cap B)$  die Wahrscheinlichkeit dafür ist, dass  $A$  und  $B$  gleichzeitig eintreten.

### 7 Definition (Exponentialverteilung):

Die Wahrscheinlichkeitsverteilung mit der Dichte

$$f_\lambda(x) := \begin{cases} 0 & \text{für } x < 0 \\ \lambda e^{-\lambda x} & \text{für } x \geq 0 \end{cases}$$

und der Verteilungsfunktion

$$F_\lambda(x) := \begin{cases} 0 & \text{für } x < 0 \\ 1 - e^{-\lambda x} & \text{für } x \geq 0 \end{cases}$$

heißt Exponentialverteilung mit Parameter  $\lambda > 0$ .

### 8 Satz:

Die Exponentialverteilung erfüllt (1).

### Beweis:

$$\begin{aligned} P(X \leq t_1 | X > t_0) &= \frac{P(\{X \leq t_1\} \cap \{X > t_0\})}{P(X > t_0)} = \frac{P(t_0 < X \leq t_1)}{P(X > t_0)} = \frac{F(t_1) - F(t_0)}{1 - F(t_0)} \\ &= \frac{1 - e^{-\lambda t_1} - (1 - e^{-\lambda t_0})}{1 - (1 - e^{-\lambda t_0})} = \frac{e^{-\lambda t_0} - e^{-\lambda t_1}}{e^{-\lambda t_0}} = \frac{1 - e^{-\lambda t_1}/e^{-\lambda t_0}}{1} \\ &= 1 - e^{-\lambda(t_1 - t_0)} = F(t_1 - t_0) = P(X \leq t_1 - t_0). \quad \blacksquare \end{aligned}$$

### 9 Bemerkung:

Es lässt sich auch die umgekehrte Richtung zeigen, also dass die Exponentialverteilung die einzige kontinuierliche Verteilung ist, die (1) erfüllt. Der Nachweis hierfür ist jedoch deutlich aufwendiger.

## 4 Modellierung eines Warteschlangensystems als Markovkette

Warteschlangensysteme mit exponentiell verteilten Zwischenankunfts- und Bedienzeiten lassen sich als Geburts- und Todesprozesse auffassen. Die Ankunft eines Kunden stellt eine Geburt dar und der Abschluss einer Bedienung einen Tod. Da Kunden auch in einem leeren System eintreffen können, handelt es sich um einen Geburts- und Todesprozess mit Einwanderung.

### 4.1 Zustandsabhängige Zwischenankunfts- und Bedienzeiten

Betrachtet werden soll ein M/M/c/K-Bediensystem, d.h. ein System mit exponentiell verteilten Zwischenankunfts- und Bedienzeiten,  $c$  parallelen Bedienern und  $K$  Warte- und Bedienplätzen. Die Ankunftsrate, d.h. die mittlere Anzahl an Ankünften pro Stunde, betrage  $\lambda$ , die mittlere Bedienrate pro Stunde  $\mu$ . Für die zustandsabhängigen Ankunfts- und Bedienraten ergeben sich damit:

$$\lambda_i := \begin{cases} \lambda & \text{für } 0 \leq i < K \\ 0 & \text{für } i = K \end{cases} \quad \text{und} \quad \mu_i := \begin{cases} i\mu & \text{für } i < c \\ c\mu & \text{für } i \geq c. \end{cases}$$

Befinden sich bereits  $K$  Kunden im System, so können keine weiteren Kunden mehr eintreffen, was zu  $\lambda_K := 0$  führt. Sind weniger als  $c$  Kunden im System, so werden diese alle gleichzeitig bedient ( $\lambda_i := i\mu$  für  $i = 0 \dots c-1$ ). Befinden sich jedoch  $c$  oder mehr Kunden im System, so sind alle Bediener ausgelastet und es werden genau  $c$  Kunden gleichzeitig bedient ( $\lambda_i := c\mu$  für  $i = c \dots K$ ).



Für den Zustand  $n$  und den Zeitschritt vom Zeitpunkt  $t$  zum Zeitpunkt  $t + \Delta t$  ergibt sich:

$$\begin{aligned}\pi_n(t + \Delta t) &= \pi_n(t)[1 - \lambda_n \Delta t - \mu_n \Delta t + o(\Delta t)] \\ &\quad + \pi_{n-1}(t)[\lambda_{n-1} \Delta t + o(\Delta t)] \\ &\quad + \pi_{n+1}(t)[\mu_{n+1} \Delta t + o(\Delta t)].\end{aligned}\tag{2}$$

Die Wahrscheinlichkeit, sich zum Zeitpunkt  $t + \Delta t$  im Zustand  $n$  ( $n = 1 \dots K - 1$ ) aufzuhalten ( $\pi_n(t + \Delta t)$ ), berechnet sich also aus der Wahrscheinlichkeit, zum Zeitpunkt  $t$  im Zustand  $n$  zu sein ( $\pi_n(t)$ ) multipliziert mit der Wahrscheinlichkeit, diesen Zustand in der Zeit  $\Delta t$  nicht zu verlassen zuzüglich der Wahrscheinlichkeiten zum Zeitpunkt  $t$  in den Zuständen  $n - 1$  bzw.  $n + 1$  gewesen zu sein ( $\pi_{n-1}(t)$  bzw.  $\pi_{n+1}(t)$ ) multipliziert mit der Wahrscheinlichkeit, dass in der Zeit  $\Delta t$  ein Kunde eingetroffen ist bzw. fertig bedient wurde.

Zieht man  $\pi_n(t)$  auf beiden Seiten von (2) ab und teilt durch  $\Delta t$ , so ergibt sich:

$$\frac{\pi_n(t + \Delta t) - \pi_n(t)}{\Delta t} = \pi_n(t) \frac{-\lambda_n \Delta t - \mu_n \Delta t + o(\Delta t)}{\Delta t} + \pi_{n-1}(t) \frac{\lambda_{n-1} \Delta t + o(\Delta t)}{\Delta t} + \pi_{n+1}(t) \frac{\mu_{n+1} \Delta t + o(\Delta t)}{\Delta t}.$$

Führt man nun den Grenzübergang  $\Delta t \rightarrow 0$  durch, so ergibt sich:

$$\pi'_n(t) = \pi_n(t)(-\lambda_n - \mu_n) + \pi_{n-1}(t)\lambda_{n-1} + \pi_{n+1}(t)\mu_{n+1}.$$

Analog ergibt sich für die Randpunkte  $n = 0$  und  $n = K$ :

$$\begin{aligned}\pi'_0(t) &= -\lambda_0 \pi_0(t) + \mu_1 \pi_1(t), \\ \pi'_K(t) &= \lambda_{K-1} \pi_{K-1}(t) - \mu_K \pi_K(t).\end{aligned}$$

Im eingeschwungenen Zustand ändern sich die Wahrscheinlichkeiten, mit denen man sich in einem bestimmten Zustand befindet, nicht mehr, d.h. es gilt  $\pi'_n(t) = 0$  für  $n = 0 \dots K$  und  $t \rightarrow \infty$ . Damit ist also folgendes Gleichungssystem zu lösen:

$$\begin{aligned}-\lambda_0 \pi_0 + \mu_1 \pi_1 &= 0, \\ \lambda_{n-1} \pi_{n-1} - (\lambda_n + \mu_n) \pi_n + \mu_{n+1} \pi_{n+1} &= 0, \quad n = 1, \dots, K - 1, \\ \lambda_{K-1} \pi_{K-1} - \mu_K \pi_K &= 0.\end{aligned}\tag{3}$$

Aus der ersten Zeile folgt sofort

$$\pi_1 = \frac{\lambda_0}{\mu_1} \pi_0\tag{4}$$

und aus der zweiten Zeile folgt für  $n = 1$

$$\pi_2 = \pi_1 \frac{\lambda_1 + \mu_1}{\mu_2} - \pi_0 \frac{\lambda_0}{\mu_2}.\tag{5}$$

Löst man (4) nach  $\pi_0$  auf und setzt dies in (5) ein, so ergibt sich:

$$\pi_2 = \pi_1 \frac{\lambda_1 + \mu_1}{\mu_2} - \pi_1 \frac{\mu_1 \lambda_0}{\lambda_0 \mu_2} = \frac{\lambda_1}{\mu_2} \pi_1,$$

was die Induktionsannahme

$$\pi_n = \frac{\lambda_{n-1}}{\mu_n} \pi_{n-1}\tag{6}$$

nahelegt. Löst man die zweite Zeile von (3) nach  $n + 1$  auf und nimmt an, dass die Induktionsannahme für  $1 \dots n$  gilt, so ergibt sich:

$$\pi_{n+1} = \pi_n \frac{\lambda_n + \mu_n}{\mu_{n+1}} - \pi_{n-1} \frac{\lambda_{n-1}}{\mu_{n+1}} = \pi_n \frac{\lambda_n + \mu_n}{\mu_{n+1}} - \pi_n \frac{\mu_n \lambda_{n-1}}{\lambda_{n-1} \mu_{n+1}} = \pi_n \frac{\lambda_n}{\mu_{n+1}}.$$

Durch rekursive Anwendung von (6) ergibt sich auch sofort:

$$\pi_n = \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0,$$

was auch für die letzte Zeile von (3), die bisher noch nicht berücksichtigt wurde, gilt.

Damit lassen sich jetzt alle Zustandswahrscheinlichkeiten  $\pi_n$ ,  $n = 1 \dots K$ , in Abhängigkeit von  $\pi_0$  ausdrücken. Da jedoch  $\sum_{n=0}^K \pi_n = 1$  gelten muss, lässt sich auch  $\pi_0$  explizit bestimmen:

$$1 = \sum_{n=0}^K \pi_n = \pi_0 + \sum_{n=1}^K \pi_n = \pi_0 + \sum_{n=1}^K \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0 = \pi_0 \left( 1 + \sum_{n=1}^K \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \right).$$

Löst man diese Gleichung nach  $\pi_0$  auf, so ergibt sich:

$$\pi_0 = \left( 1 + \sum_{n=1}^K \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \right)^{-1}.$$

Insgesamt ergeben sich damit für die stationären Zustandswahrscheinlichkeiten:

$$\pi_n = \begin{cases} \left( 1 + \sum_{m=1}^K \prod_{i=1}^m \frac{\lambda_{i-1}}{\mu_i} \right)^{-1} & \text{für } n = 0, \\ \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0 & \text{für } n = 1 \dots K. \end{cases}$$

### 5.3 Bestimmung der Verteilung der Anzahl an Kunden im System im stationären Fall

Wählt man  $K = \infty$ , d.h. lässt man die Warteraumbegrenzung weg, so lässt sich auch die Verteilung für die Anzahl an Kunden im System relativ einfach berechnen. (Ohne die Annahme  $K = \infty$  würde die  $\Gamma$ -Funktion in dem Ergebnis auftreten.)

Die Wahrscheinlichkeitsverteilung  $P(W \leq t)$  für die Wartezeiten der Kunden im stationären Fall lässt sich zunächst mit Hilfe des Satzes von der totalen Wahrscheinlichkeit wie folgt beschreiben:

$$P(W \leq t) = \sum_{n=0}^{\infty} P(W \leq t | N = n) \cdot \pi_n$$

Ist die Anzahl der Kunden im System kleiner oder gleich  $c - 1$ , so ist mindestens ein Bediener frei und ein eintreffender Kunde kann sofort bedient werden, d.h. für den Fall  $n \leq c - 1$  gilt  $P(W \leq t | N = n) = 1$  für jedes  $t \geq 0$ . Außerdem ist  $\sum_{n=0}^{c-1} \pi_n$  gerade die Wahrscheinlichkeit dafür, dass ein neu eintreffender Kunde nicht warten muss, also  $P(W = 0)$ . Damit folgt:

$$\begin{aligned} P(W \leq t) &= \sum_{n=0}^{\infty} P(W \leq t | N = n) \cdot \pi_n = \sum_{n=0}^{c-1} 1 \cdot \pi_n + \sum_{n=c}^{\infty} P(W \leq t | N = n) \cdot \pi_n \\ &= P(W = 0) + \sum_{n=c}^{\infty} P(n - c + 1 \text{ Kunden bedient } \leq t | N = n) \cdot \pi_n. \end{aligned}$$

Im Fall  $n \geq c$  sind alle Bediener ausgelastet, d.h. die Bedienrate pro Zeiteinheit lautet  $c\mu$ . Der Takt, mit dem die Warteschlange abgearbeitet wird, entspricht dem Minimum von  $c$  stochastisch unabhängigen und mit Parameter  $\mu > 0$  exponentiell verteilten Zufallsgrößen. (Die Bedienzeiten der wartenden

Kunden sind exponentiell verteilt; die Restbedienzeiten der Kunden in Bedienung genügen aufgrund der Gedächtnislosigkeit der Exponentialverteilung ebenfalls der Exponentialverteilung.) Damit entspricht die Dauer für die Bedienung von  $n - c + 1$  Kunden einer Erlang( $n - c + 1, c\mu$ ) verteilten Zufallsgröße. Und es folgt:

$$P(W \leq t) = P(W = 0) + \sum_{n=c}^{\infty} \int_0^t f_{\text{Erlang}(n-c+1, c\mu)}(x) dx \cdot \pi_n.$$

Mit  $f_{\text{Erlang}(k, \lambda)}(x) := \frac{\lambda^k}{(k-1)!} x^{k-1} e^{-\lambda x}$  für  $x \geq 0$  folgt weiter:

$$P(W \leq t) = P(W = 0) + \sum_{n=c}^{\infty} \int_0^t \frac{(c\mu)^{n-c+1}}{(n-c)!} x^{n-c} e^{-c\mu x} \pi_0 \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} dx.$$

Für  $n \geq c$  war

$$\prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} = \prod_{i=1}^c \frac{\lambda}{i\mu} \cdot \prod_{i=c+1}^n \frac{\lambda}{c\mu} = \frac{\lambda^c}{c! \mu^c} \cdot \left( \frac{\lambda}{c\mu} \right)^{n-c}.$$

Damit:

$$\begin{aligned} P(W \leq t) &= P(W = 0) + \pi_0 \frac{\lambda^c}{\mu^c c!} \sum_{n=c}^{\infty} \int_0^t \frac{(c\mu)^{n-c+1}}{(n-c)!} x^{n-c} e^{-c\mu x} \left( \frac{\lambda}{c\mu} \right)^{n-c} dx \\ &= P(W = 0) + \pi_0 \frac{\lambda^c}{\mu^c c!} \int_0^t c\mu \cdot e^{-c\mu x} \underbrace{\sum_{n=c}^{\infty} \frac{(c\mu x \frac{\lambda}{\mu c})^{n-c}}{(n-c)!}}_{=\sum_{n=0}^{\infty} \frac{(\lambda x)^n}{n!} = e^{\lambda x}} dx \\ &= P(W = 0) + \pi_0 \frac{\lambda^c}{\mu^c c!} \int_0^t c\mu \cdot e^{-c\mu x} \cdot e^{\lambda x} dx \\ &= P(W = 0) + \pi_0 \frac{\lambda^c}{\mu^{c-1} (c-1)!} \left[ -\frac{1}{c\mu - \lambda} e^{-(c\mu - \lambda)x} \right]_{x=0}^{x=t} \\ &= P(W = 0) + \pi_0 \frac{\lambda^c}{\mu^{c-1} (c-1)! (c\mu - \lambda)} \left( 1 - e^{-(c\mu - \lambda)t} \right). \end{aligned}$$

Nun gilt aber:

$$\begin{aligned} P(W > 0) &= P(N \geq c) = \sum_{n=c}^{\infty} \pi_n = \sum_{n=c}^{\infty} \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0 = \sum_{n=c}^{\infty} \underbrace{\prod_{i=1}^c \frac{\lambda}{i\mu}}_{=\frac{\lambda^c}{c! \mu^c}} \cdot \underbrace{\prod_{i=c+1}^n \frac{\lambda}{c\mu}}_{=(\frac{\lambda}{c\mu})^{n-c}} \pi_0 \\ &= \frac{\lambda^c}{c! \mu^c} \pi_0 \sum_{n=1}^{\infty} \left( \frac{\lambda}{c\mu} \right)^n = \frac{\lambda^c}{c! \mu^c} \cdot \frac{1}{1 - \frac{\lambda}{c\mu}} \pi_0 = \pi_0 \frac{\lambda^c}{(c-1)! \mu^{c-1} (c\mu - \lambda)}. \end{aligned}$$

Und damit schließlich:

$$\begin{aligned} P(W \leq t) &= P(W = 0) + P(W > 0) \left( 1 - e^{-(c\mu - \lambda)t} \right) = 1 - P(W > 0) + P(W > 0) \left( 1 - e^{-(c\mu - \lambda)t} \right) \\ &= 1 - P(W > 0) \cdot e^{-(c\mu - \lambda)t}. \end{aligned}$$

## 6 Berechnung weiterer Kenngrößen

Mit Hilfe der stationären Verteilung lassen sich eine Reihe von Kenngrößen berechnen. Dafür wird insbesondere die Formel von LITTLE (siehe Satz 5) verwendet.

- Mittlere Anzahl an Kunden im System:

$$E[N] = \sum_{n=0}^K n \cdot \pi_n = \sum_{n=0}^K n \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0.$$

- Mittlere Verweilzeit (Summe aus Wartezeit und Bedienzeit):

$$E[V] = \frac{1}{\lambda} E[N] = \frac{1}{\lambda} \sum_{n=0}^K n \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0.$$

- Mittlere Anzahl an Kunden in der Warteschlange:

$$E[N_Q] = \sum_{n=c+1}^K (n-c) \cdot \pi_n = \sum_{n=c+1}^K (n-c) \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0.$$

- Mittlere Wartezeit der Kunden:

$$E[W] = \frac{1}{\lambda} E[N_Q] = \frac{1}{\lambda} \sum_{n=c+1}^K (n-c) \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i} \pi_0.$$

- Besetztwahrscheinlichkeit, d.h. die Wahrscheinlichkeit, dass ein neu eintreffender Kunde abgewiesen wird:

$$P(B) = \pi_K = \prod_{i=1}^K \frac{\lambda_{i-1}}{\mu_i} \pi_0.$$

- Wahrscheinlichkeit, dass ein Kunde warten muss (Blockierwahrscheinlichkeit):

$$P(W > 0) = \pi_0 \frac{\lambda^c}{(c-1)! \mu^{c-1} (c\mu - \lambda)}.$$

- Service-Level zum  $t$ -Sekunden-Niveau (Wahrscheinlichkeit, dass ein Kunde höchstens  $t$ -Sekunden warten muss):

$$P(W \leq t) = 1 - P(W > 0) \cdot e^{-(c\mu - \lambda)t}.$$

Häufig soll ein System so ausgelegt werden, dass 80% der Kunden innerhalb von weniger als 20 Sekunden bedient werden oder sogar 90% in weniger als 10 Sekunden.