

Single Queuing Systems

M/M/1 queuing system

- arrival process is a Poisson process (or the inter-arrival time is exponentially distributed)
- service process is also a Poisson process (or the service time is exponentially distributed)

A counting process $\{N(t), t \geq 0\}$

representing the
of events
that have
occurred
up to time t

advantage: a mathematically tractable model with solutions applicable to a wide variety of situations.

Poisson process with an average arrival rate λ : λ is the proportionality constant



$$\Pr(\text{exactly 1 arrival in } [t, t+\Delta t]) = \lambda\Delta t$$

$$\Pr(\text{no arrivals in } [t, t+\Delta t]) = 1 - \lambda\Delta t$$

$$\Pr \left\{ \begin{array}{l} \text{1 event} \\ \text{occurs at } t + \Delta t \end{array} \middle| \begin{array}{l} \text{event does not} \\ \text{occur at } t \end{array} \right\}$$

$$= 1 - e^{-\lambda\Delta t}$$

$$= 1 - \left\{ 1 + (-\lambda\Delta t) + \frac{(-\lambda\Delta t)^2}{2!} + \dots \right\}$$

$$\approx \lambda\Delta t$$

Analogy:

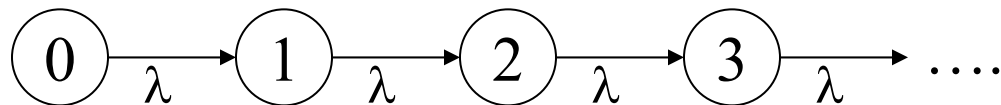
Coin flipping: results of coin flips are independent

Arrivals are also independent

Let $P_n(t) \equiv P(\text{\# of arrivals} = n \text{ at time } t)$

$P_{ij}(\Delta t) \equiv$ the prob. of going from i arrivals to j arrivals in a time interval of Δt seconds

$$\begin{aligned} \therefore P_n(t + \Delta t) &= P_n(t) \overbrace{P_{n,n}(\Delta t)}^{1-\lambda\Delta t} + P_{n-1}(t) \underbrace{P_{n-1,n}(\Delta t)}_{\equiv \lambda\Delta t} \\ P_0(t + \Delta t) &= P_0(t) \underbrace{P_{0,0}(\Delta t)}_{1-\lambda\Delta t} \\ \frac{P_n(t + \Delta t) - P_n(t)}{\Delta t} &= -\lambda P_n(t) + \lambda P_{n-1}(t) \\ \frac{P_0(t + \Delta t) - P_0(t)}{\Delta t} &= -\lambda P_0(t) \end{aligned}$$



Let $\Delta t \rightarrow 0$ then

$$\frac{dP_n(t)}{dt} = -\lambda P_n(t) + \lambda P_{n-1}(t) \text{ ——— } \textcircled{1}$$

$$\frac{dP_0(t)}{dt} = -\lambda P_0(t) \text{ ——— } \textcircled{2} \xrightarrow{\text{solution}} P_0(t) = e^{-\lambda t}$$

From $\textcircled{1}$

$$\frac{dP_1(t)}{dt} = -\lambda P_1(t) + \lambda e^{-\lambda t} \quad \therefore P_1 = \lambda t e^{-\lambda t}$$

$$\frac{dP_2(t)}{dt} = -\lambda P_2(t) + \lambda^2 t e^{-\lambda t} \quad \therefore P_2 = \frac{\lambda^2 t^2}{2} e^{-\lambda t}$$

Continuing, by induction,

$$P_n(t) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}$$

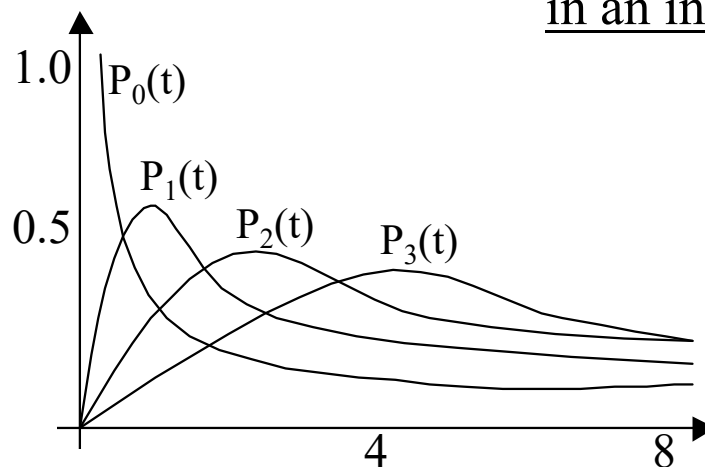
— Poisson distribution

meaning: prob. of n arrivals
in an interval of t seconds

$$\sum_{n=0}^{\infty} P_n(t) = 1$$

$$\sum_{n=0}^{\infty} \frac{(\lambda t)^n}{n!} e^{-\lambda t} = 1$$

$$\sum_{n=0}^{\infty} \frac{X^n}{n!} = e^X$$



The Poisson distribution

Ex: $\lambda = 100$ arrivals/min., what is the prob. of no arrivals in 5 sec.?

$$P_0(5 \text{ sec.}) = e^{-100 \cdot \left(\frac{1}{12}\right)} = 0.00024$$

* the mean & the variance of the Poisson dist. are both equal to λt .

mean:

$$\overline{n(t)} = \sum_{n=1}^{\infty} n P_n(t) = \sum_{n=1}^{\infty} n \cdot \frac{(\lambda t)^n e^{-\lambda t}}{n!} = e^{-\lambda t} \cdot \sum_{n=1}^{\infty} \frac{(\lambda t)^n}{(n-1)!}$$

derivation is based on

$$= e^{-\lambda t} \cdot \sum_{n=0}^{\infty} \frac{(\lambda t)^{n+1}}{n!} = e^{-\lambda t} \cdot \lambda t \cdot \sum_{n=0}^{\infty} \frac{(\lambda t)^n}{n!} = e^{-\lambda t} \cdot \lambda t \cdot e^{\lambda t} = \lambda t$$

$$e^{\lambda t} = \sum_{n=0}^{\infty} \frac{(\lambda t)^n}{n!}$$

variance: $[\sigma_n(t)]^2 = \sum_{n=0}^{\infty} n^2 P_n(t) - [\overline{n(t)}]^2 = \lambda t$

The inter-arrival time T ← a random variable

Inter-arrival time cumulative dist. function(t) ← cdf (t)

= P (time between arrivals $\leq t$)

= 1 - P (time between arrivals $> t$) \because no arrivals in a time interval of t = $P_0(t)$

= 1 - $\underbrace{P_0(t)}_{P_0(t)}$

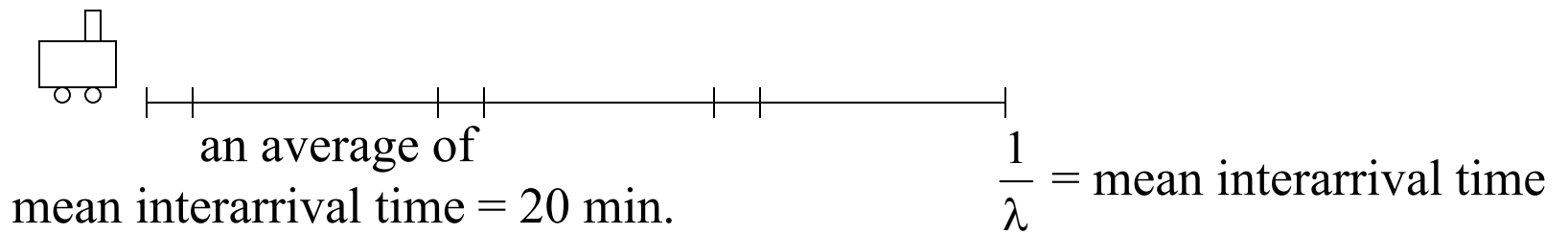
= $1 - e^{-\lambda t}$

\therefore inter-arrival time density (t) = $\frac{d(1 - e^{-\lambda t})}{dt} = \lambda e^{-\lambda t}$ ← pdf (t)

\therefore T is an exponentially distributed r.v.

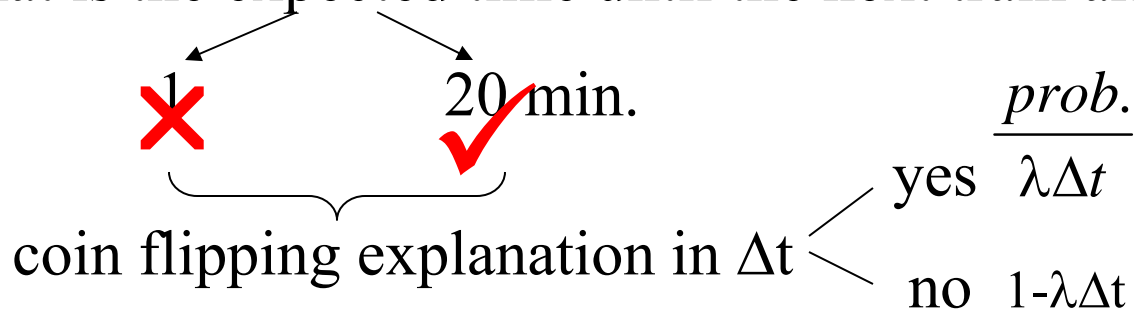
\therefore T has a memory-less property

e.g.,



The last train arrived 19 minutes ago.

What is the expected time until the next train arrives?



Memory-less property \equiv Markov property

$$P (T > t_0+t | T > t_0) = P (T > t)$$

* Definition: a Markov chain is a Markov process with a discrete state space.

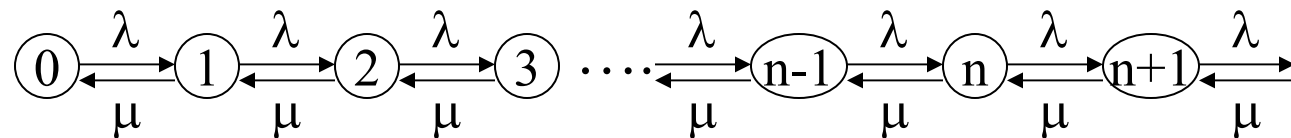
M/M/1

$$P_n(t + \Delta t) = P_n(t) \overbrace{P_{n,n}(\Delta t)}^{(1-\lambda\Delta t)(1-\mu\Delta t)} + P_{n-1}(t) \overbrace{P_{n-1,n}(\Delta t)}^{\lambda\Delta t} + P_{n+1}(t) \overbrace{P_{n+1,n}(\Delta t)}^{\mu\Delta t}$$

$$P_0(t + \Delta t) = P_0(t) \overbrace{P_{0,0}(\Delta t)}^{1-\lambda\Delta t} + P_1(t) \overbrace{P_{1,0}(\Delta t)}^{\mu\Delta t}$$

$$\frac{dP_n(t)}{dt} = -(\lambda + \mu)P_n(t) + \lambda P_{n-1}(t) + \mu P_{n+1}(t)$$

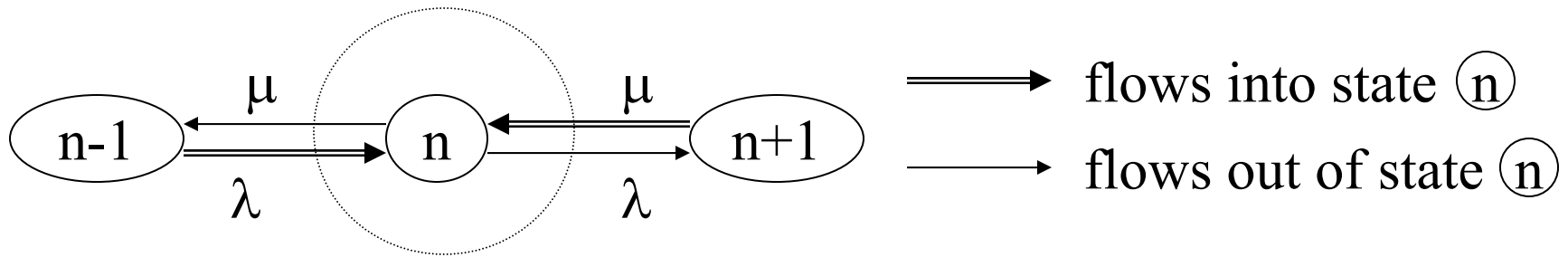
$$\frac{dP_0(t)}{dt} = -\lambda P_0(t) + \mu P_1(t)$$



Probability flux (or flow):

(probability of a state) * (transition rate originating from the state)

physical meaning: # of times per second the event corresponding to the transition occurs.



$$\frac{dP_n(t)}{dt} = \underbrace{-(\lambda + \mu)P_n(t)}_{\text{flows out of state n}} + \underbrace{\lambda P_{n-1}(t) + \mu P_{n+1}(t)}_{\text{flows into state n}}$$

Study:

— transient behavior $\frac{dP_n(t)}{dt} \neq 0$

— equilibrium behavior $\frac{dP_n(t)}{dt} = \frac{dP_{n-1}(t)}{dt} = \dots = \frac{dP_0(t)}{dt} = 0$

this yields

Global balance equations: a set of linear equations for $t \rightarrow \infty$

Equilibrium state probabilities

conservation of probability:

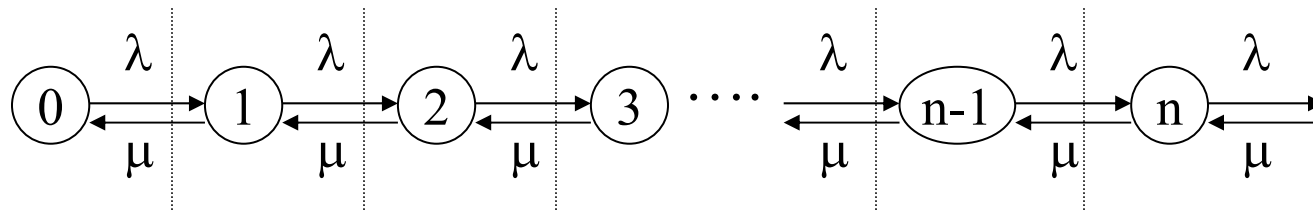
$$\sum_{i=0}^{\infty} P_i(\infty) = 1 \longrightarrow \boxed{\text{normalization equation}}$$

Use local balance equations to solve the global balance equations

1. Local satisfies global
2. Local allows us to relate P_n with a reference state, e.g., P_0

Definition of local balance:

“the probability flow into a state due to an arrival to a queue equals the probability flow out of the same state due to a departure from the same queue”



$$\begin{aligned}
P_0: & P_0 \\
P_1: & P_0\lambda = P_1\mu \quad \Rightarrow P_1 = \left(\frac{\lambda}{\mu}\right)P_0 \\
P_2: & P_1\lambda = P_2\mu \quad \Rightarrow P_2 = \left(\frac{\lambda}{\mu}\right)P_1 \\
& \vdots \\
& P_{n-1}\lambda = P_n\mu \quad \Rightarrow P_n = \left(\frac{\lambda}{\mu}\right)P_{n-1} \\
\hline
& P_n = \left(\frac{\lambda}{\mu}\right)^n P_0
\end{aligned}$$

applying the normalization equation

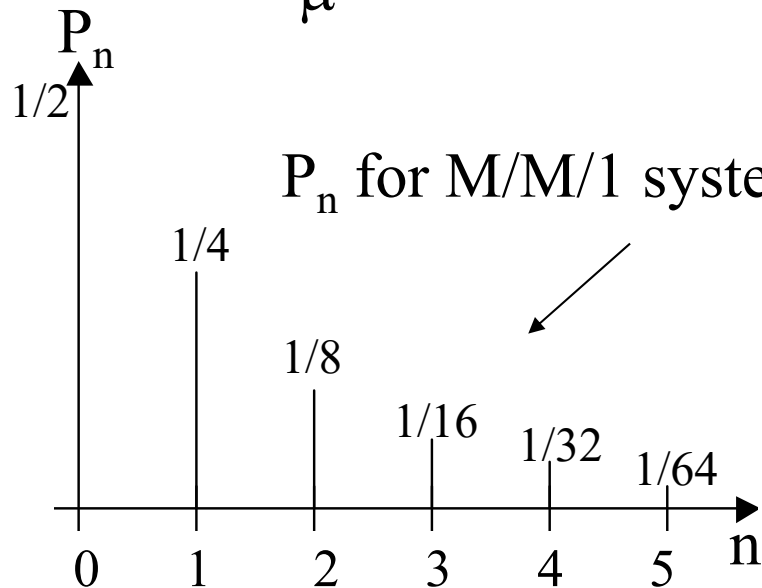
$$\therefore P_0 \sum_{n=0}^{\infty} \left(\frac{\lambda}{\mu}\right)^n = 1 \quad \sum_{n=0}^{\infty} x^n = \frac{1}{1-x} \quad 0 \leq x < 1$$

$$P_0 = \frac{1}{\sum_{n=0}^{\infty} \left(\frac{\lambda}{\mu}\right)^n} = 1 - \left(\frac{\lambda}{\mu}\right) \quad \text{if } 0 \leq \frac{\lambda}{\mu} < 1$$

$$\therefore P_n = \left(\frac{\lambda}{\mu}\right)^n \left(1 - \frac{\lambda}{\mu}\right)$$

Utilization: prob. that an M/M/1 queuing system is nonempty

$$\text{Let } \frac{\lambda}{\mu} = \rho, \quad P_n = \rho^n (1 - \rho)$$



* for a lightly loaded system, there are usually less than 4 customers in the system.

	(light) * $\rho=0.1$	$\rho=0.5$	(heavy) $\rho=0.9$
$P(0 \leq n \leq 3)$	0.9999	0.9375	0.3439
$P(4 \leq n \leq 7)$	10^{-4}	0.0586	0.2256
$P(8 \leq n \leq 11)$	10^{-8}	0.0366	0.1480
$P(n \geq 12)$	10^{-12}	0.000249	0.2825

check $\rho = 1 - P_0$?

$$\rho = 1 - P_0 = 1 - \left(1 - \frac{\lambda}{\mu}\right) = \frac{\lambda}{\mu}$$

* $\rho \leq 1$ otherwise $\lambda > \mu$ and the queuing system would no longer be in equilibrium \rightarrow i.e., unstable.

Q1: throughput?

$$\begin{aligned} x &= \sum_{n=1}^{\infty} P_n * \mu \\ &= \mu \cdot \sum_{n=1}^{\infty} P_n \\ &= \mu(1 - P_0) \\ &= \mu \cdot \frac{\lambda}{\mu} = \lambda \end{aligned}$$

because when there is no customer, there is no contribution to throughput.

Q2: Average # of customers in the queuing system?

$$\begin{aligned} \bar{n} &= \sum_{n=0}^{\infty} n P_n = \sum_{n=0}^{\infty} n \cdot \rho^n (1 - \rho) = (1 - \rho) \sum_{n=0}^{\infty} n \cdot \rho^n = (1 - \rho) \cdot \rho \sum_{n=0}^{\infty} n \cdot \rho^{n-1} \\ &= (1 - \rho) \cdot \rho \frac{d}{d\rho} \sum_{n=0}^{\infty} \rho^n = (1 - \rho) \cdot \rho \cdot \frac{d}{d\rho} \left(\frac{1}{1 - \rho}\right) = (1 - \rho) \cdot \rho \cdot \frac{1}{(1 - \rho)^2} \end{aligned}$$

$$= \frac{\rho}{1 - \rho}$$

$$\bar{n} = \frac{\rho}{1 - \rho}$$

e.g.,

$$\lambda = \frac{1}{2} \mu$$

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

$$= \frac{\frac{1}{2} \mu}{\mu} = \frac{1}{2}$$

$$\bar{n} = \frac{\frac{1}{2}}{\frac{1}{2}} = 1$$

e.g.,

$$\lambda = \frac{2}{3} \mu$$

$$\therefore \rho = \frac{2}{3}$$

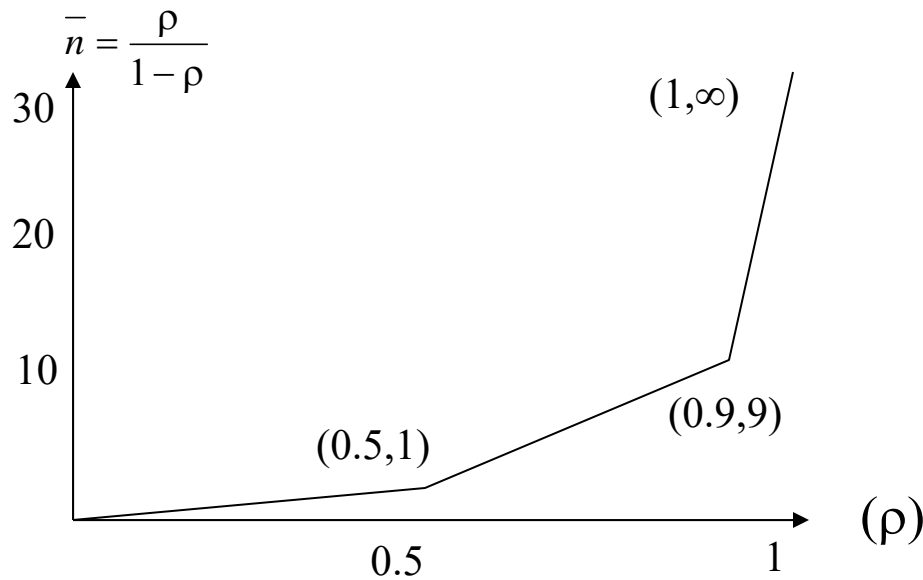
$$\bar{n} = \frac{\frac{2}{3}}{\frac{1}{3}} = 2$$

Let R be the mean response time per customer

Q3: R ?

since $\bar{n} = \lambda R$ by little's law (to be discussed later)

$$R = \frac{\bar{n}}{\lambda} = \frac{\rho}{(1-\rho)\lambda} = \frac{1/\mu}{(1-\frac{\lambda}{\mu})} \quad \longrightarrow \quad \text{When } \rho=1 \text{ system is unstable}$$

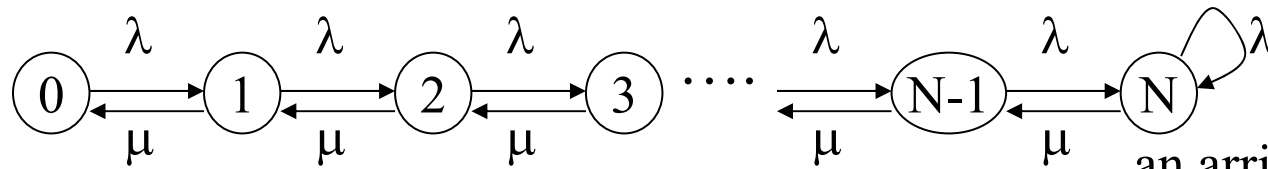


M/M/1: average # of customers \bar{n}
as a function of ρ

M/M/1 service time
waiting time

$$\begin{aligned} R &= (1 + \bar{n}) \cdot D \\ &= (1 + \bar{n}) \cdot \frac{1}{\mu} \\ &= \left(1 + \frac{\rho}{1-\rho} \right) \cdot \frac{1}{\mu} \\ &= \frac{1}{1-\rho} \cdot \left(\frac{1}{\mu} \right) \end{aligned}$$

M/M/1/N Queuing system: the finite buffer case



an arriving customer is “lost” or “turned away” when there are already N customers in the system.

Following the previous derivation for M/M/1/∞,

$$P_n = \left(\frac{\lambda}{\mu}\right)^n P_0$$

&

$$P_0 = \frac{1}{\sum_{n=0}^N \left(\frac{\lambda}{\mu}\right)^n} = \frac{1}{\frac{1 - \left(\frac{\lambda}{\mu}\right)^{N+1}}{1 - \left(\frac{\lambda}{\mu}\right)}} = \frac{1 - \left(\frac{\lambda}{\mu}\right)}{1 - \left(\frac{\lambda}{\mu}\right)^{N+1}}$$

no restriction on the range of $\frac{\lambda}{\mu}$

$$\therefore P_n = \frac{1 - \left(\frac{\lambda}{\mu}\right)}{1 - \left(\frac{\lambda}{\mu}\right)^{N+1}} \cdot \left(\frac{\lambda}{\mu}\right)^n$$

$$\begin{array}{ll} N \nearrow & P_N \searrow \\ \lambda / \mu \nearrow & P_N \nearrow \end{array}$$

Q1: the prob. that the queuing system is full? P_N

Q2: how fast are customers lost? $P_N \times \lambda$

	$\frac{\lambda}{\mu}$	$P_N = P_5$ (blocking probability)
when N=5	0.1	$9 \cdot 10^{-6}$
	0.5	0.016
	0.75	0.072
	1.00	0.166
	2.00	0.508
	5.00	0.800

applying L'Hopital's rule

$$\lim_{x \rightarrow 1} \frac{x^5 - x^6}{1 - x^6} = \lim_{x \rightarrow 1} \frac{5x^4 - 6x^5}{-6x^5} = 0.166$$

Q3: population?

$$\bar{n} = \sum_{n=0}^N n \cdot P_n$$

Q4: throughput?

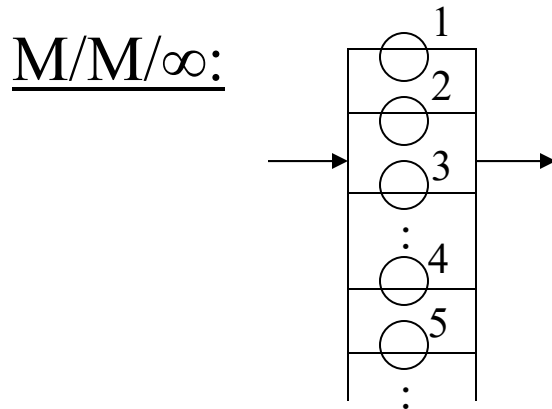
$$x = \sum_{n=1}^N \mu \cdot P_n = \mu \cdot \sum_{n=1}^N P_n = \mu(1 - P_0) = \mu\rho < \lambda$$

Q5: Utilization?

$$\rho = 1 - P_0 = 1 - \frac{1 - \frac{\lambda}{\mu}}{1 - \left(\frac{\lambda}{\mu}\right)^{N+1}} < \frac{\lambda}{\mu}$$

ρ is utilization

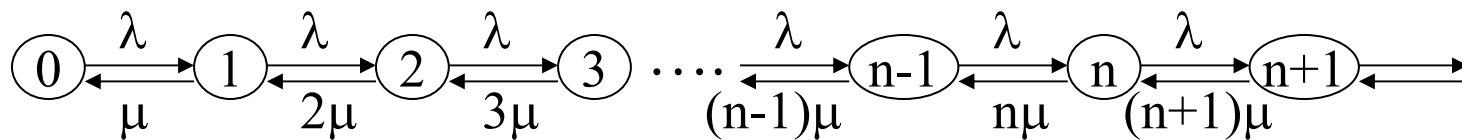
Variations of M/M/1



infinite # of servers

- Q1: throughput? λ
 - Q2: response time? $1/\mu$
 - Q3: population? λ/μ
- by Little's law

$$\bar{n} = \sum_{n=1}^{\infty} n \cdot P_n = \sum_{n=1}^{\infty} n \frac{\left(\lambda \cdot \frac{1}{\mu}\right)^n e^{-\lambda \cdot \frac{1}{\mu}}}{n!} = \frac{\lambda}{\mu}$$



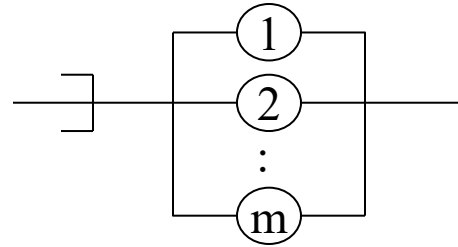
$$P_n = \left(\prod_{i=1}^n \frac{\lambda}{i\mu} \right) \cdot P_0 = \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n P_0$$

$$\therefore P_0 = \frac{1}{\sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n} = \frac{1}{e^{+(\frac{\lambda}{\mu})}} = e^{-\frac{\lambda}{\mu}}$$

a Poisson process
with mean $\bar{n} = \frac{\lambda}{\mu}$

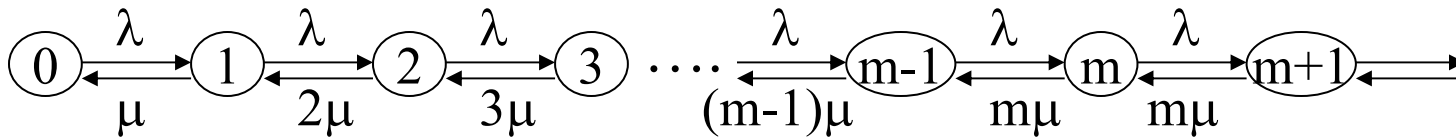
$$P_n = \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n e^{-\lambda/\mu}$$

M/M/m



m servers

e.g., a system with m processors



Solution:

$$P_n = \frac{1}{m! m^{n-m}} \left(\frac{\lambda}{\mu} \right)^n \cdot P_0$$

can be obtained by
consider 2 cases
separately

$$\mu_n = \begin{cases} n\mu & 0 \leq n \leq m \\ m\mu & n \geq m \end{cases}$$

where $P_0 = \left[1 + \sum_{n=1}^{m-1} \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n + \frac{1}{m!} \left(\frac{\lambda}{\mu} \right)^m \left(\frac{1}{1-\rho} \right) \right]^{-1}$

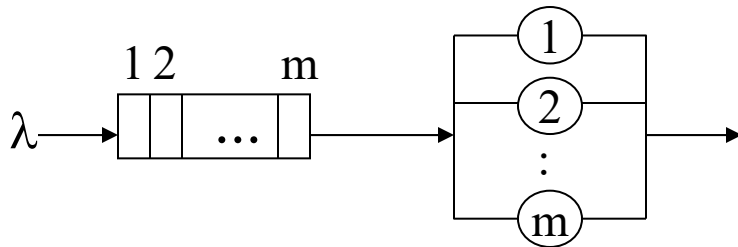
where $\rho = \frac{\lambda}{m\mu}$

Q1: what is the probability that all servers are busy? Ans: $\sum_{n=m}^{\infty} P_n$

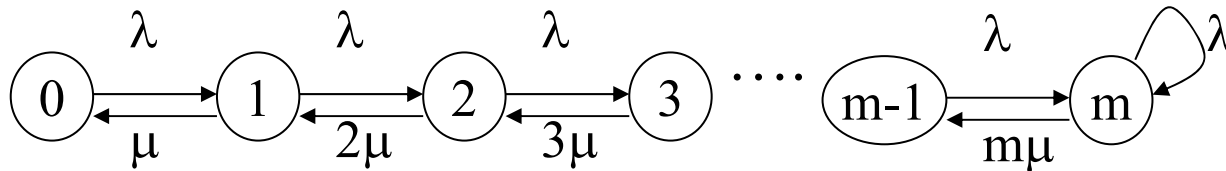
Q2: throughput? Ans: λ

Q3: response time?

M/M/m/m



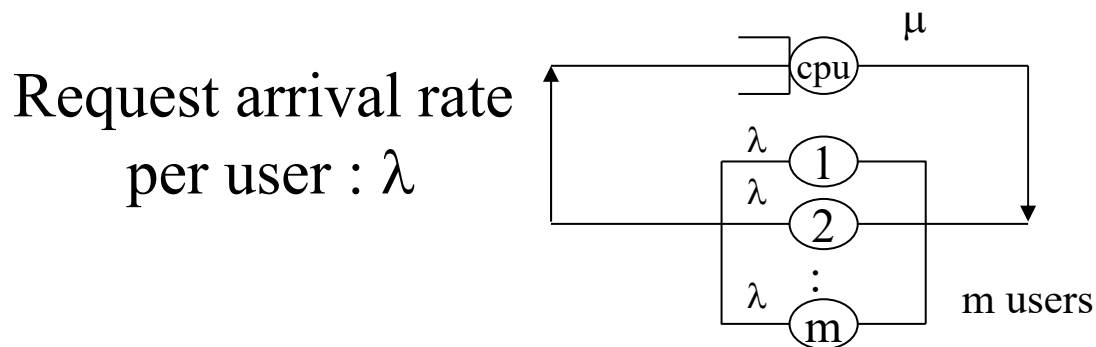
m servers with a single queue having a buffer space of m (when all servers are busy, a customer walks away), e.g., a telephone switching system.



$$P_n = \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n P_0 \quad \therefore P_0 = \frac{1}{\sum_{n=0}^m \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n} \quad \& \quad P_n = \frac{\frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^n}{\sum_{i=0}^m \frac{1}{n!} \left(\frac{\lambda}{\mu} \right)^i}$$

- Q1. Prob. that all m servers are busy (e.g., in a telephone switch company)? P_m ← The expression for P_m is called Erlang's B formula.
- Q2. Mean # of calls turned away per time unit? $P_m \times \lambda$

A Client-Server System

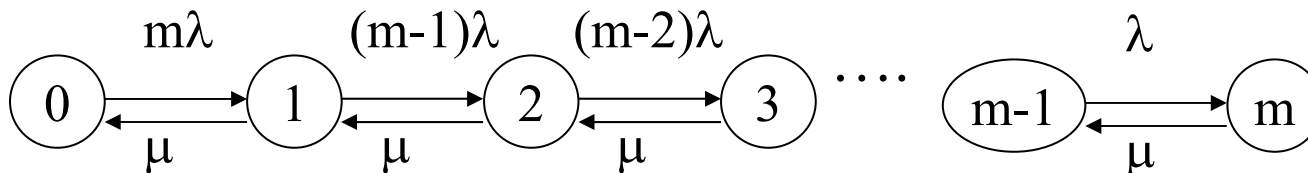


Response time: the time spent by a user at the system between submitting the request & the return of the response

State Description: one state component representation

n : a number representing the # of users in the server system

\therefore # of users still thinking (i.e., not issuing requests) = $m - n$



Recall
in M/M/1

$$P_n = \left(\frac{\lambda}{\mu}\right)^n P_0$$

$$\therefore P_n = \left(\prod_{i=1}^n \frac{\lambda(m-i+1)}{\mu}\right) P_0$$

$$\text{or } P_n = \left(\frac{\lambda}{\mu}\right)^n \frac{m!}{(m-n)!} P_0$$

$$\text{where } P_0 = \frac{1}{\sum_{n=0}^m \left[\left(\frac{\lambda}{\mu}\right)^n \frac{m!}{(m-n)!} \right]}$$

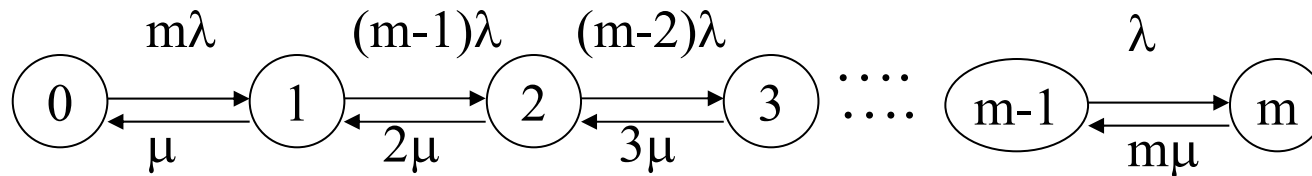
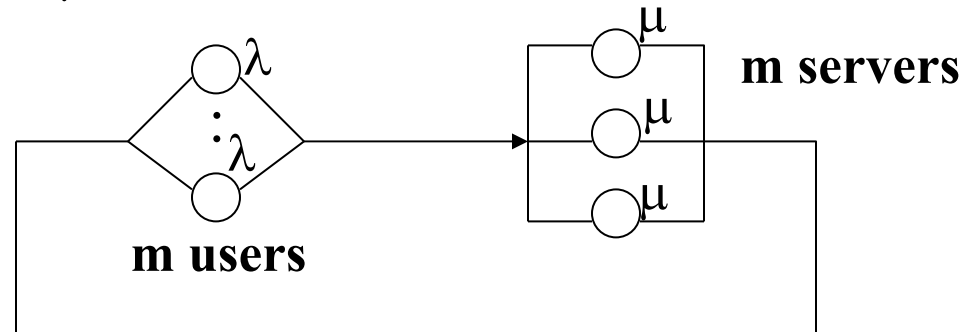
Q1: Avg. # of users in the server system? $\bar{n} = \sum_{n=1}^m n \cdot P_n$

Q2: Avg. # of users still thinking (not issuing requests)? $m - \bar{n}$

Q3: System throughput? $x = \sum_{n=1}^m \mu \cdot P_n = (1 - P_0) \cdot \mu$

Q4: Response time per user?

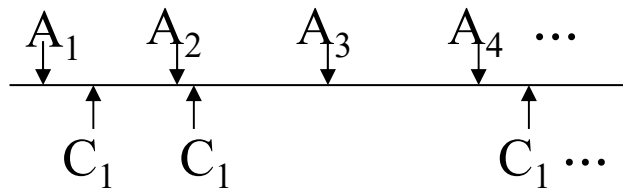
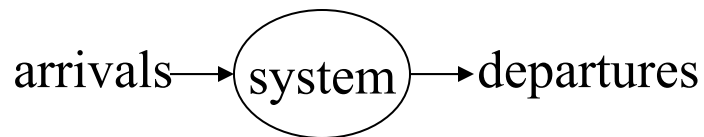
What happens if the server system has m servers, each with a service rate of μ ?



Q1: Throughput? $x = \sum_{n=1}^m (P_n * n\mu)$

Q2: Response time?

Fundamental Laws: algebraic relationships among performance measurement quantities.



$\lambda = \text{arrival rate} = A/T$ e.g., $\frac{100 \text{ arrivals}}{1 \text{ hr}}$

$C = \# \text{ of completions}$

$x = C/T$ throughput

$B = \text{total system busy time}$

$D = B/C$ average service time per request

$\rho = B/T$ utilization of the system

mathematically

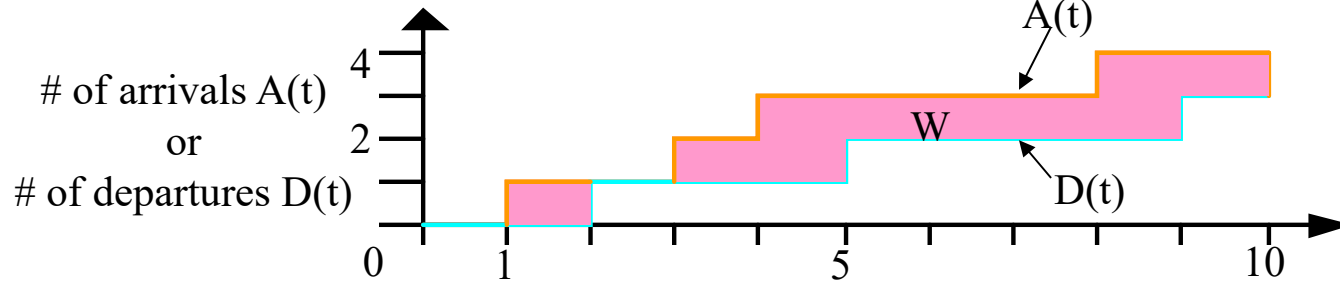
$$\frac{B}{T} = \frac{C}{T} * \frac{B}{C} \quad \boxed{\rho = x * D}$$

utilization law

Little's law

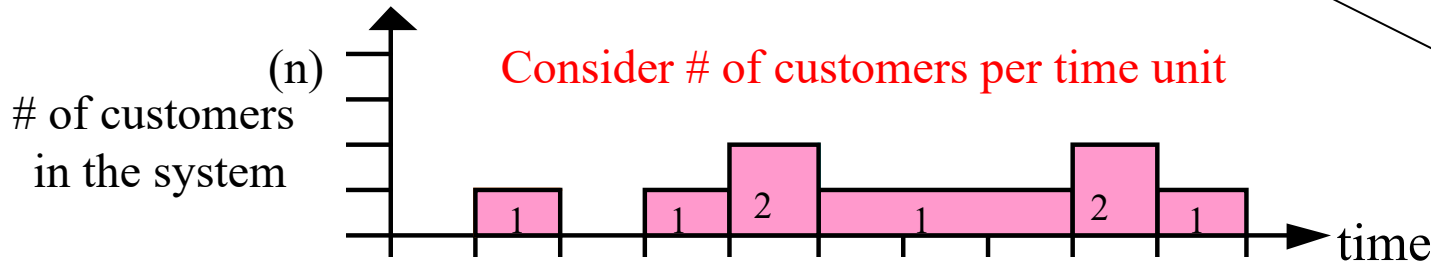
Consider response time per customer

A	D
1	2
3	5
4	9
8	10



the shaded region
 $\bar{n} = \frac{W}{T}$
 observation period

Consider # of customers per time unit



* A meaning of W is the total time spent by all customers in the system. $\therefore R = W/C$

* Another meaning of W is the total population accumulated (in queue & in service) over T time units. $\therefore \bar{n} = \frac{W}{T}$

Algebraically $\bar{n} = \frac{W}{T} = \frac{W}{C} * \overbrace{\frac{C}{T}}^x = R * x$ $\therefore \bar{n} = R x$