

CS 147:
Computer Systems Performance Analysis
Introduction to Queueing Theory

2015-06-15 CS147

CS 147:
Computer Systems Performance Analysis
Introduction to Queueing Theory

Overview

Introduction and Terminology

Poisson Distributions

Fundamental Results

Stability

Little's Law

M/M/*

M/M/1

M/M/m

M/M/m/B

More General Queues

2015-06-15 CS147

└ Overview

Overview

Introduction and Terminology
Poisson Distributions

Fundamental Results
Stability
Little's Law

M/M/*

M/M/1

M/M/m

M/M/m/B

More General Queues

What is a Queueing System?

- ▶ A *queueing system* is any system in which things arrive, hang around for a while, and leave
- ▶ Examples
 - ▶ A bank
 - ▶ A freeway
 - ▶ A (computer) network
 - ▶ A beehive
- ▶ The things that arrive and leave are *customers* or *jobs*
- ▶ Customers leave after receiving *service*
- ▶ Most queueing systems have (surprise!) a *queue* that can store (delay) customers awaiting service

2015-06-15

CS147

└ Introduction and Terminology

└ What is a Queueing System?

What is a Queueing System?

- A *queueing system* is any system in which things arrive, hang around for a while, and leave
- Examples
 - A bank
 - A freeway
 - A (computer) network
 - A beehive
- The things that arrive and leave are *customers* or *jobs*
- Customers leave after receiving *service*
- Most queueing systems have (surprise!) a *queue* that can store (delay) customers awaiting service

Parameters of a Queueing System

Arrival Process Injects customers into system

- ▶ Usually statistical
- ▶ Convenient to specify in terms of *interarrival time* distribution
- ▶ Most common is *Poisson* arrivals

Service Time Also statistical

Number of Servers Often 1

System Capacity Equals number of servers plus queue capacity.
Often assumed infinite for convenience

Population Maximum number of customers. Often assumed infinite

Service Discipline How next customer is chosen for service.
Often FCFS or priority

2015-06-15

CS147

└ Introduction and Terminology

└ Parameters of a Queueing System

Parameters of a Queueing System

Arrival Process Injects customers into system

- Usually statistical
- Convenient to specify in terms of *interarrival time* distribution
- Most common is *Poisson* arrivals

Service Time Also statistical

Number of Servers Often 1

System Capacity Equals number of servers plus queue capacity.
Often assumed infinite for convenience

Population Maximum number of customers. Often assumed infinite

Service Discipline How next customer is chosen for service.
Often FCFS or priority

Arrival and Service Distributions

- ▶ Customer arrivals are random variables
 - ▶ Next disk request from many processes
 - ▶ Next packet hitting Google
 - ▶ Next call to Chipotle
- ▶ Same is true for service times
- ▶ What distribution describes it?
 - ▶ May be complicated (fractal, Zipf)
 - ▶ We often use Poisson for tractability

2015-06-15

CS147

└ Introduction and Terminology

└ Arrival and Service Distributions

Arrival and Service Distributions

- Customer arrivals are random variables
 - Next disk request from many processes
 - Next packet hitting Google
 - Next call to Chipotle
- Same is true for service times
- What distribution describes it?
 - May be complicated (fractal, Zipf)
 - We often use Poisson for tractability

The Poisson Distribution

- ▶ Probability of exactly k arrivals in $(0, t)$ is $P_k(t) = (\lambda t)^k e^{-\lambda t} / k!$
 - ▶ λ is *arrival rate* parameter
- ▶ More useful formulation is *Poisson arrival distribution*:
 - ▶ PDF $A(t) = P[\text{next arrival takes time} \leq t] = 1 - e^{-\lambda t}$
 - ▶ pdf $a(t) = \lambda e^{-\lambda t}$
 - ▶ Also known as *exponential* or *memoryless* distribution
 - ▶ Mean = standard deviation = λ
- ▶ Poisson distribution is *memoryless*
 - ▶ Assume $P[\text{arrival within 1 second}]$ at time $t_0 = x$
 - ▶ Then $P[\text{arrival within 1 second}]$ at time $t_1 > t_0$ is also x
 - ▶ I.e., no memory that time has passed
 - ▶ Often true in real world
 - ▶ E.g., when I go to Von's doesn't affect when you go

2015-06-15

CS147

└ Introduction and Terminology

└ Poisson Distributions

└ The Poisson Distribution

The Poisson Distribution

- Probability of exactly k arrivals in $(0, t)$ is $P_k(t) = (\lambda t)^k e^{-\lambda t} / k!$
- λ is *arrival rate* parameter
- More useful formulation is *Poisson arrival distribution*:
 - PDF $A(t) = P[\text{next arrival takes time} \leq t] = 1 - e^{-\lambda t}$
 - pdf $a(t) = \lambda e^{-\lambda t}$
 - Also known as *exponential* or *memoryless* distribution
 - Mean = standard deviation = λ
- Poisson distribution is *memoryless*
 - Assume $P[\text{arrival within 1 second}]$ at time $t_0 = x$
 - Then $P[\text{arrival within 1 second}]$ at time $t_1 > t_0$ is also x
 - I.e., no memory that time has passed
 - Often true in real world
 - E.g., when I go to Von's doesn't affect when you go

Splitting and Merging Poisson Processes

- ▶ Merging streams of Poisson events (e.g., arrivals) is Poisson

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

- ▶ Splitting a Poisson stream randomly gives Poisson streams; if stream i has probability p_i , then

$$\lambda_i = p_i \lambda$$

2015-06-15

CS147

└ Introduction and Terminology

└ Poisson Distributions

└ Splitting and Merging Poisson Processes

Splitting and Merging Poisson Processes

- Merging streams of Poisson events (e.g., arrivals) is Poisson

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

- Splitting a Poisson stream randomly gives Poisson streams; if stream i has probability p_i , then

$$\lambda_i = p_i \lambda$$

Kendall's Notation

$A/S/m/B/K/D$ defines a (single) queueing system compactly:

- A Denotes arrival distribution, as follows:
 - M Exponential (M emoryless)
 - E_k Erlang with parameter k
 - D Deterministic
 - G Completely general (very hard to analyze!)
- S Service distribution, same as arrival
- m Number of servers
- B System capacity; ∞ if omitted
- K Population size; ∞ if omitted
- D Service discipline, FCFS if omitted

2015-06-15

CS147

└ Introduction and Terminology

└ Poisson Distributions

└ Kendall's Notation

Kendall's Notation

$A/S/m/B/K/D$ defines a (single) queueing system compactly:

- A Denotes arrival distribution, as follows:
 - M Exponential (M emoryless)
 - E_k Erlang with parameter k
 - D Deterministic
 - G Completely general (very hard to analyze!)
- S Service distribution, same as arrival
- m Number of servers
- B System capacity; ∞ if omitted
- K Population size; ∞ if omitted
- D Service discipline, FCFS if omitted

Examples of Kendall's Notation

D/D/1 Arrivals on clock tick, fixed service times, one server

M/M/m Memoryless arrivals, memoryless service, multiple servers (good model of a bank)

M/M/m/m Customers go away rather than wait in line

G/G/1 Modern disk drive

2015-06-15

CS147

└ Introduction and Terminology

└ Poisson Distributions

└ Examples of Kendall's Notation

Examples of Kendall's Notation

D/D/1 Arrivals on clock tick, fixed service times, one server
M/M/m Memoryless arrivals, memoryless service, multiple servers (good model of a bank)
M/M/m/m Customers go away rather than wait in line
G/G/1 Modern disk drive

Common Variables

- τ Interarrival time. Usually varies per customer, e.g., τ_1, τ_2, \dots
- λ Mean arrival rate: $1/\bar{\tau}$
- s_i Service time for job i , sometimes called x_i
- μ Mean service rate per server, $1/\bar{s}$
- ρ Traffic intensity or system load = $\lambda/m\mu$. **This is the most important parameter in most queueing systems**
- w_j Waiting time, or time in queue: interval between arrival and beginning of service
- r_j Response time = $w_j + s_j$

2015-06-15 CS147
 └ Introduction and Terminology
 └ Poisson Distributions
 └ Common Variables

Common Variables

- Interarrival time. Usually varies per customer, e.g., τ_1, τ_2, \dots
- Mean arrival rate: $1/\bar{\tau}$
- Service time for job i , sometimes called x_i
- Mean service rate per server: $1/\bar{s}$
- Traffic intensity or system load = $\lambda/m\mu$. **This is the most important parameter in most queueing systems**
- Waiting time, or time in queue: interval between arrival and beginning of service
- Response time = $w_j + s_j$

Stability

- ▶ A system is stable iff $\lambda < m\mu \equiv \rho < 1$
- ▶ Otherwise, system can't keep up and queue grows to ∞
- ▶ Exception: in D/D/m, $\rho = 1$ is OK

2015-06-15
CS147
└ Fundamental Results
└└ Stability
└└└ Stability

Stability

- A system is stable iff $\lambda < m\mu \equiv \rho < 1$
- Otherwise, system can't keep up and queue grows to ∞
- Exception: in D/D/m, $\rho = 1$ is OK

Little's Law

- ▶ Let n = Number of jobs in system
- ▶ Then $n = \lambda \bar{r}$
- ▶ Likewise, if n_q = Number of jobs in queue, then $n_q = \lambda \bar{w}$
- ▶ True regardless of distributions, queueing disciplines, etc., as long as system is in equilibrium
- ▶ May seem obvious:
 - ▶ If ten people are ahead of you in line, and each takes about 1 minute for service, you're going to be stuck there for 10 minutes
- ▶ Not proved until 1961
- ▶ Often useful for calculating queue lengths:
 - ▶ Packet takes 2s to arrive, you're sending 100 pps
 - ⇒ Mean queue length = $100 \text{ pkt/s} \times 2\text{s} = 200 \text{ pkts}$

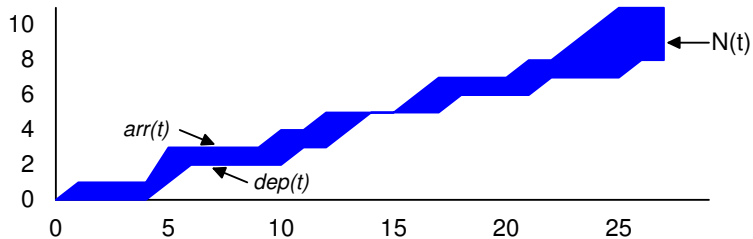
2015-06-15
 CS147
 └ Fundamental Results
 └ └ Little's Law
 └ └ └ Little's Law

Little's Law

- Let n = Number of jobs in system
- Then $n = \lambda \bar{r}$
- Likewise, if n_q = Number of jobs in queue, then $n_q = \lambda \bar{w}$
- True regardless of distributions, queueing disciplines, etc., as long as system is in equilibrium
- May seem obvious:
 - If ten people are ahead of you in line, and each takes about 1 minute for service, you're going to be stuck there for 10 minutes
- Not proved until 1961
- Often useful for calculating queue lengths:
 - Packet takes 2s to arrive, you're sending 100 pps
 - Mean queue length = $100 \text{ pkt/s} \times 2\text{s} = 200 \text{ pkts}$

Deriving Little's Law

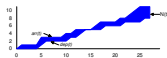
- ▶ Define $arr(t) = \#$ of arrivals in interval $(0, t)$
- ▶ Define $dep(t) = \#$ of departures in interval $(0, t)$
- ▶ Clearly, $N(t) = \#$ in system at time $t = arr(t) - dep(t)$
- ▶ Area between curves = $spent(t) =$ total time spent in system by all customers (measured in customer-seconds)



2015-06-15
 CS147
 └ Fundamental Results
 └ Little's Law
 └ Deriving Little's Law

Deriving Little's Law

- Define $arr(t) = \#$ of arrivals in interval $(0, t)$
- Define $dep(t) = \#$ of departures in interval $(0, t)$
- Clearly, $N(t) = \#$ in system at time $t = arr(t) - dep(t)$
- Area between curves = $spent(t) =$ total time spent in system by all customers (measured in customer-seconds)



Deriving Little's Law (continued)

- ▶ Define average arrival rate during interval t , in customers/second, as $\lambda_t = arr(t)/t$
- ▶ Define T_t as system time/customer, averaged over all customers in $(0, t)$
 - ▶ Since $spent(t) =$ accumulated customer-seconds, divide by arrivals up to that point to get $T_t = spent(t)/arr(t)$
- ▶ Mean tasks in system over $(0, t)$ is accumulated customer-seconds divided by seconds:
 $Mean-tasks_t = spent(t)/t$
- ▶ Above three equations give us:

$$\begin{aligned}
 Mean-tasks_t &= spent(t)/t \\
 &= T_t arr(t)/t \\
 &= \lambda_t T_t
 \end{aligned}$$

2015-06-15

CS147

Fundamental Results

Little's Law

Deriving Little's Law (continued)

Deriving Little's Law (continued)

- Define average arrival rate during interval t , in customers/second, as $\lambda_t = arr(t)/t$
- Define T_t as system time/customer, averaged over all customers in $(0, t)$
 - Since $spent(t) =$ accumulated customer-seconds, divide by arrivals up to that point to get $T_t = spent(t)/arr(t)$
- Mean tasks in system over $(0, t)$ is accumulated customer-seconds divided by seconds:
Mean-tasks _{t} = $spent(t)/t$
- Above three equations give us:

$$\begin{aligned}
 Mean-tasks_t &= spent(t)/t \\
 &= T_t arr(t)/t \\
 &= \lambda_t T_t
 \end{aligned}$$

Deriving Little's Law (continued)

- ▶ We've shown that $Mean\text{-}tasks_t = \lambda_t T_t$
- ▶ Assuming limits of λ_t and T_t exist, limit of $mean\text{-}tasks_t$ also exists and gives Little's result:

Mean tasks in system = arrival rate \times mean time in system

2015-06-15 CS147
└ Fundamental Results
└ Little's Law
└ Deriving Little's Law (continued)

Deriving Little's Law (continued)

- We've shown that $Mean\text{-}tasks_t = \lambda_t T_t$
- Assuming limits of λ_t and T_t exist, limit of $mean\text{-}tasks_t$ also exists and gives Little's result:

Mean tasks in system = arrival rate \times mean time in system

The M/M/1 Queue

- ▶ Remember this one if you don't remember anything else
- ▶ Assumptions are sometimes realistic, sometimes not
 - ▶ Never infinite customers or capacity
 - ▶ Service times aren't truly Poisson
 - ▶ Interarrival times more likely to be Poisson
- ▶ Still provides surprisingly good analysis
- ▶ M/M/1's characteristics are clue to many other queues
- ▶ Primary results (in equilibrium):
 - ▶ Mean number in system $\bar{n} = \rho/(1 - \rho)$
 - ▶ Mean time in system

$$\bar{r} = (1/\mu)/(1 - \rho) = 1/\mu(1 - \rho) = \bar{s}/(1 - \rho)$$

2015-06-15
 CS147
 └ M/M/*
 └ M/M/1
 └ The M/M/1 Queue

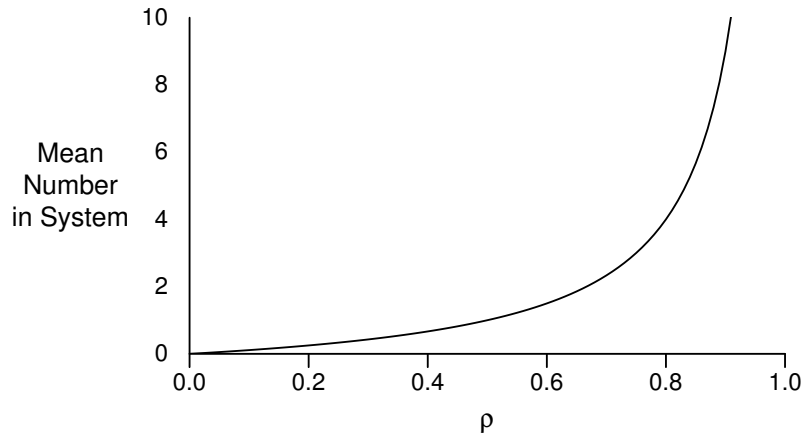
The M/M/1 Queue

- Remember this one if you don't remember anything else
- Assumptions are sometimes realistic, sometimes not
 - Never infinite customers or capacity
 - Service times aren't truly Poisson
 - Interarrival times more likely to be Poisson
- Still provides surprisingly good analysis
- M/M/1's characteristics are clue to many other queues
- Primary results (in equilibrium):
 - Mean number in system $\bar{n} = \rho/(1 - \rho)$
 - Mean time in system

$$\bar{r} = (1/\mu)/(1 - \rho) = 1/\mu(1 - \rho) = \bar{s}/(1 - \rho)$$

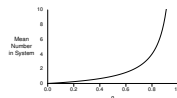
Nearly all useful results in queueing theory apply only to systems in equilibrium.

The Nastiness of High Load



2015-06-15
CS147
└ M/M/*
└ M/M/1
└ The Nastiness of High Load

The Nastiness of High Load



The system breaks down completely at $\rho > 0.95$.

The reason for the breakdown is variance: at high load, a burst fills the queue and it takes a long time to drain, giving plenty of time for another burst to arrive.

More M/M/1 Results

- ▶ Variance is $\rho/(1 - \rho)^2$, so standard deviation is $\sqrt{\rho}/(1 - \rho)$
- ▶ q -percentile of time in system is $\bar{r} \ln[100/(100 - q)]$
 - ▶ 90th percentile is $2.3\bar{r}$
- ▶ Mean *waiting* time is $\bar{w} = \frac{1}{\mu} \frac{\rho}{1 - \rho}$
- ▶ q -percentile of waiting time is

$$\max \left(0, \frac{\bar{w}}{\rho} \ln[100\rho/(100 - q)] \right)$$

- ▶ Mean jobs served in a busy period: $1/(1 - \rho)$
- ▶ Probability of n jobs in system $\rho_n = (1 - \rho)\rho^n$
- ▶ Probability of $> n$ jobs in system: ρ^n

2015-06-15

CS147

└ M/M/*

└ M/M/1

└ More M/M/1 Results

More M/M/1 Results

- Variance is $\rho/(1 - \rho)^2$, so standard deviation is $\sqrt{\rho}/(1 - \rho)$
- q -percentile of time in system is $\bar{r} \ln[100/(100 - q)]$
 - 90th percentile is $2.3\bar{r}$
- Mean waiting time is $\bar{w} = \frac{1}{\mu} \frac{\rho}{1 - \rho}$
- q -percentile of waiting time is

$$\max \left(0, \frac{\bar{w}}{\rho} \ln[100\rho/(100 - q)] \right)$$
- Mean jobs served in a busy period: $1/(1 - \rho)$
- Probability of n jobs in system $\rho_n = (1 - \rho)\rho^n$
- Probability of $> n$ jobs in system: ρ^n

M/M/1 Example

- ▶ Web server gets 1200 requests/hour w/ Poisson arrivals
- ▶ Typical request takes 1s to serve
- ▶ $\rho = 1200/3600 = 0.33$
- ▶ Mean requests in service = $0.33/0.67 = 0.5$
- ▶ Mean response time $\bar{r} = (1/1)/(1 - 0.33) = 1.5s$
- ▶ 90th percentile response time = 3.4s

2015-06-15

CS147

└ M/M/*

└ M/M/1

└ M/M/1 Example

MM/1 Example

- Web server gets 1200 requests/hour w/ Poisson arrivals
- Typical request takes 1s to serve
- $\rho = 1200/3600 = 0.33$
- Mean requests in service = $0.33/0.67 = 0.5$
- Mean response time $\bar{r} = (1/1)/(1 - 0.33) = 1.5s$
- 90th percentile response time = 3.4s

This slide has animations.

M/M/1 Example

- ▶ Web server gets 1200 requests/hour w/ Poisson arrivals
- ▶ Typical request takes 1s to serve
- ▶ $\rho = 1200/3600 = 0.33$
- ▶ Mean requests in service = $0.33/0.67 = 0.5$
- ▶ Mean response time $\bar{r} = (1/1)/(1 - 0.33) = 1.5s$
- ▶ 90th percentile response time = 3.4s
- ▶ But if Slashdot hits. . .

2015-06-15
CS147
└─ M/M/*
 └─ M/M/1
 └─ M/M/1 Example

M/M/1 Example

- Web server gets 1200 requests/hour w/ Poisson arrivals
- Typical request takes 1s to serve
- $\rho = 1200/3600 = 0.33$
- Mean requests in service = $0.33/0.67 = 0.5$
- Mean response time $\bar{r} = (1/1)/(1 - 0.33) = 1.5s$
- 90th percentile response time = 3.4s
- But if Slashdot hits. . .

This slide has animations.

M/M/1 Example (cont'd)

- ▶ Suppose Slashdot raises request rate to 3500/hr
- ▶ Now $\rho = 3500/3600 = 0.972$
- ▶ Mean requests in service = $0.972/(1 - 0.972) = 34.7$
- ▶ $\bar{r} = 1/0.028 = 35.7$ seconds
- ▶ 90th percentile response time = 82.8s

2015-06-15

CS147

└ M/M/*

└ M/M/1

└ M/M/1 Example (cont'd)

M/M/1 Example (cont'd)

- Suppose Slashdot raises request rate to 3500/hr
- Now $\rho = 3500/3600 = 0.972$
- Mean requests in service = $0.972/(1 - 0.972) = 34.7$
- $\bar{r} = 1/0.028 = 35.7$ seconds
- 90th percentile response time = 82.8s

This slide has animations.

M/M/1 Example (cont'd)

- ▶ Suppose Slashdot raises request rate to 3500/hr
- ▶ Now $\rho = 3500/3600 = 0.972$
- ▶ Mean requests in service = $0.972/(1 - 0.972) = 34.7$
- ▶ $\bar{r} = 1/0.028 = 35.7$ seconds
- ▶ 90th percentile response time = 82.8s
- ▶ And don't even *think* about 4000 requests/hr

2015-06-15

CS147

└ M/M/*

└ M/M/1

└ M/M/1 Example (cont'd)

M/M/1 Example (cont'd)

- Suppose Slashdot raises request rate to 3500/hr
- Now $\rho = 3500/3600 = 0.972$
- Mean requests in service = $0.972/(1 - 0.972) = 34.7$
- $\bar{r} = 1/0.028 = 35.7$ seconds
- 90th percentile response time = 82.8s
- And don't even think about 4000 requests/hr

This slide has animations.

M/M/m

- ▶ Multiple servers, one queue
- ▶ $\rho = \lambda / (m\mu)$
- ▶ We'll need probability of empty system:

$$p_0 = \frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}}$$

- ▶ Probability of queueing:

$$\rho = P(\geq m \text{ jobs}) = \frac{(m\rho)^m}{m!(1-\rho)} p_0$$

2015-06-15

CS147
 └ M/M/*
 └ M/M/m
 └ M/M/m

MM/m

- Multiple servers, one queue
- $\rho = \lambda / (m\mu)$
- We'll need probability of empty system:

$$p_0 = \frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}}$$

- Probability of queueing:

$$\rho = P(\geq m \text{ jobs}) = \frac{(m\rho)^m}{m!(1-\rho)} p_0$$

For $m = 1$, $\rho = \rho$

M/M/m (cont'd)

- ▶ Mean jobs in system: $\bar{n} = m\rho + \rho\varrho/(1 - \rho)$
- ▶ Mean time in system:

$$\bar{r} = \frac{1}{\mu} \left(1 + \frac{\varrho}{m(1 - \rho)} \right)$$

- ▶ Mean waiting time: $\bar{w} = \varrho/[m\mu(1 - \rho)]$
- ▶ q -percentile of waiting time:

$$\max \left(0, \frac{\bar{w}}{\varrho} \ln \frac{100\varrho}{100 - q} \right)$$

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ M/M/m (cont'd)

M/M/m (cont'd)

- Mean jobs in system: $\bar{n} = m\rho + \rho\varrho/(1 - \rho)$
- Mean time in system:

$$\bar{r} = \frac{1}{\mu} \left(1 + \frac{\varrho}{m(1 - \rho)} \right)$$

- Mean waiting time: $\bar{w} = \varrho/[m\mu(1 - \rho)]$
- q -percentile of waiting time:

$$\max \left(0, \frac{\bar{w}}{\varrho} \ln \frac{100\varrho}{100 - q} \right)$$

$m \times M/M/1$ vs. $M/M/m$

- ▶ For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
 - ▶ But ρ is unchanged
- ▶ $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- ▶ $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- ▶ $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- ▶ $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- ▶ $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ $m \times M/M/1$ vs. $M/M/m$ $m \times M/M/1$ vs. $M/M/m$

- For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
 - But ρ is unchanged
- $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

This slide has animations.

$m \times M/M/1$ vs. $M/M/m$

- ▶ For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
 - ▶ But ρ is unchanged
- ▶ $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- ▶ $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- ▶ $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- ▶ $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- ▶ $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- ▶ $\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} > \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ $m \times M/M/1$ vs. $M/M/m$ $m \times M/M/1$ vs. $M/M/m$

- For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
 - But ρ is unchanged
- $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} > \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

This slide has animations.

$m \times M/M/1$ vs. $M/M/m$

- ▶ For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
 - ▶ But ρ is unchanged
- ▶ $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- ▶ $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- ▶ $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- ▶ $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- ▶ $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- ▶ $\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} > \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ $m \times M/M/1$ vs. $M/M/m$ $m \times M/M/1$ vs. $M/M/m$

- For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
 - But ρ is unchanged
- $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} > \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

This slide has animations.

$m \times M/M/1$ vs. $M/M/m$

- ▶ For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$

- ▶ But ρ is unchanged

- ▶ $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$

- ▶ $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$

- ▶ $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$

- ▶ $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

- ▶ $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

- ▶ $\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} > \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ $m \times M/M/1$ vs. $M/M/m$ $m \times M/M/1$ vs. $M/M/m$

- For m separate $M/M/1$ queues, each queue sees arrival rate of $\lambda_{M/M/1} = \lambda/m$
- But ρ is unchanged
- $\bar{r}_{m \times M/M/1} = \frac{1}{\mu} \left(\frac{1}{1-\rho} \right) \stackrel{?}{>} \bar{r}_{M/M/m} = \frac{1}{\mu} \left(1 + \frac{\rho}{m(1-\rho)} \right)$
- $1 \stackrel{?}{>} 1 - \rho + \frac{\rho}{m}$
- $\rho \stackrel{?}{>} \rho_0 \frac{(m\rho)^m}{m!m(1-\rho)}$
- $1 \stackrel{?}{>} \rho_0 \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $1 \stackrel{?}{>} \left(\frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} \right) \frac{(m\rho)^{m-1}}{m!(1-\rho)}$
- $\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!} > \frac{(m\rho)^{m-1}}{m!(1-\rho)}$

This slide has animations.

Running Some Numbers

- ▶ Assume 5 servers, $\rho = 0.5, \mu = 1$
- ▶ Then $\bar{r}_{m \times M/M/1} = 1/(1 - \rho) = 2$
- ▶ $\varrho = \frac{(m\rho)^m}{m!(1-\rho)} \rho_0 = \frac{(2.5)^5}{5!(0.5)} \rho_0 = \frac{97.7}{60} \rho_0 = 1.63\rho_0$
- ▶ $\rho_0 = \frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} = \frac{1}{1.63+1+\frac{2.5^1}{1}+\frac{2.5^2}{2}+\frac{2.5^3}{3!}+\frac{2.5^4}{4!}}$
- ▶ $\rho_0 = \frac{1}{1.63+1+2.5+3.13+2.60+1.63} = \frac{1}{12.49} = 0.08$
- ▶ So $\varrho = 1.63(0.08) = 0.13$
- ▶ And $\bar{r}_{m/M/m} = 1 + \frac{\varrho}{m(1-\rho)} = 1 + \frac{0.13}{5(1-0.5)} = 1 + \frac{0.13}{2.5} = 1.05$
- ▶ In terms of previous slide's inequality,
 $\frac{97.7}{60} + 1 + 2.5 + 3.13 + 2.60 + 1.63 = 12.49 > \frac{2.5^4}{5!(0.5)} = \frac{39.1}{60} = 0.65$

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ Running Some Numbers

Running Some Numbers

- Assume 5 servers, $\rho = 0.5, \mu = 1$
- Then $\bar{r}_{m \times M/M/1} = 1/(1 - \rho) = 2$
- $\varrho = \frac{(m\rho)^m}{m!(1-\rho)} \rho_0 = \frac{97.7}{60} \rho_0 = 1.63\rho_0$
- $\rho_0 = \frac{1}{\frac{(m\rho)^m}{m!(1-\rho)} + \sum_{k=0}^{m-1} \frac{(m\rho)^k}{k!}} = \frac{1}{1.63+1+\frac{2.5^1}{1}+\frac{2.5^2}{2}+\frac{2.5^3}{3!}+\frac{2.5^4}{4!}}$
- $\rho_0 = \frac{1}{1.63+1+2.5+3.13+2.60+1.63} = \frac{1}{12.49} = 0.08$
- So $\varrho = 1.63(0.08) = 0.13$
- And $\bar{r}_{m/M/m} = 1 + \frac{\varrho}{m(1-\rho)} = 1 + \frac{0.13}{5(1-0.5)} = 1 + \frac{0.13}{2.5} = 1.05$
- In terms of previous slide's inequality,
 $\frac{97.7}{60} + 1 + 2.5 + 3.13 + 2.60 + 1.63 = 12.49 > \frac{2.5^4}{5!(0.5)} = \frac{39.1}{60} = 0.65$

$m \times M/M/1$ vs. $M/M/m$ (cont'd)

- ▶ A similar result holds for variance
- ▶ Conclusion: single queue, multiple server is **always** better than one queue per server
- ▶ Question 1: When is this false? (hint: multiple cores)
- ▶ Question 2: Why do so many movie theaters have multiple lines for popcorn?

2015-06-15

CS147

└ M/M/*

└ M/M/m

└ $m \times M/M/1$ vs. $M/M/m$ (cont'd) $m \times M/M/1$ vs. $M/M/m$ (cont'd)

- A similar result holds for variance
- Conclusion: single queue, multiple server is **always** better than one queue per server
- Question 1: When is this false? (hint: multiple cores)
- Question 2: Why do so many movie theaters have multiple lines for popcorn?

M/M/m/B

- ▶ Real systems have finite capacity
- ▶ Previous analysis applies only under light loads (relative to capacity)
- ▶ Considering limit has several effects:
 - ▶ Lost jobs (obviously)
 - ▶ Loss rate ρ_B becomes important parameter
 - ▶ Mean response time drops compared to $M/M/m/\infty$ (Why?)

2015-06-15
CS147
└─ M/M/*
 └─ M/M/m/B
 └─ M/M/m/B

MM/m/B

- Real systems have finite capacity
- Previous analysis applies only under light loads (relative to capacity)
- Considering limit has several effects:
 - Lost jobs (obviously)
 - Loss rate ρ_B becomes important parameter
 - Mean response time drops compared to $M/M/m/\infty$ (Why?)

Extending the Results

- ▶ Unsurprisingly, generality equates to (mathematical) complexity
- ▶ Many special cases have been analyzed (e.g., Erlang distributions)
- ▶ Little's Law always applies
- ▶ Important cases:
 - ▶ M/G/1
 - ▶ M/D/1
 - ▶ G/G/m (but mostly intractable)

2015-06-15

CS147

└ More General Queues

└ Extending the Results

Extending the Results

- Unsurprisingly, generality equates to (mathematical) complexity
- Many special cases have been analyzed (e.g., Erlang distributions)
- Little's Law always applies
- Important cases:
 - M/G/1
 - M/D/1
 - G/G/m (but mostly intractable)