

# Queuing models

Stochastic throughput

**Cathy Wu**

1.041/1.200 Transportation: Foundations and Methods

# Readings

1. Larson, Richard C. and Amedeo R. Odoni. **Urban Operations Research**. Prentice-Hall (1981). Chapter 4: Queueing Theory. [URL](#).

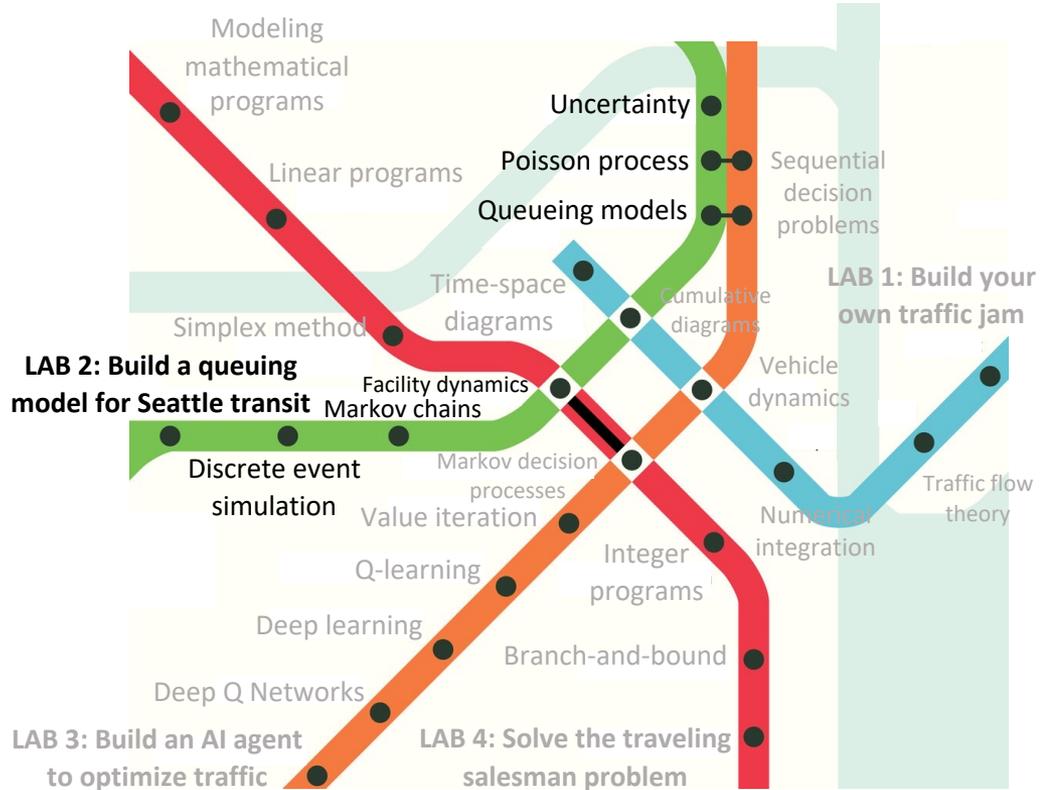
# Unit 2: Queuing systems



Unit 2

Modeling

Stochastic



# Outline

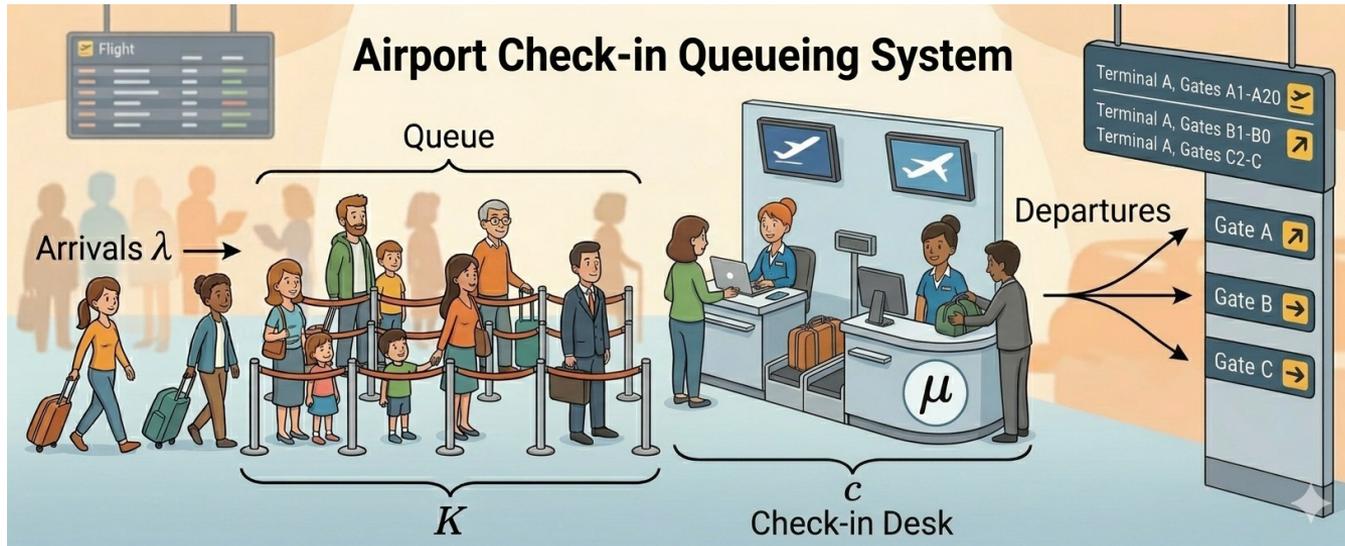
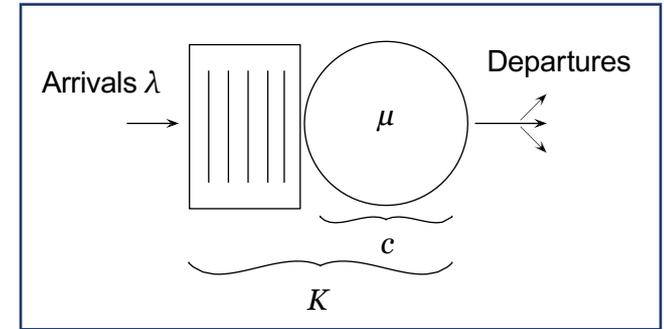
1. Fundamental queueing models
2. Stationary analysis and Little's law
3. M/M/1: Detailed analysis
4. More queues
5. Airport security checkpoints problem

# Outline

1. **Fundamental queueing models**
2. Stationary analysis and Little's law
3. M/M/1: Detailed analysis
4. More queues
5. Airport security checkpoints problem

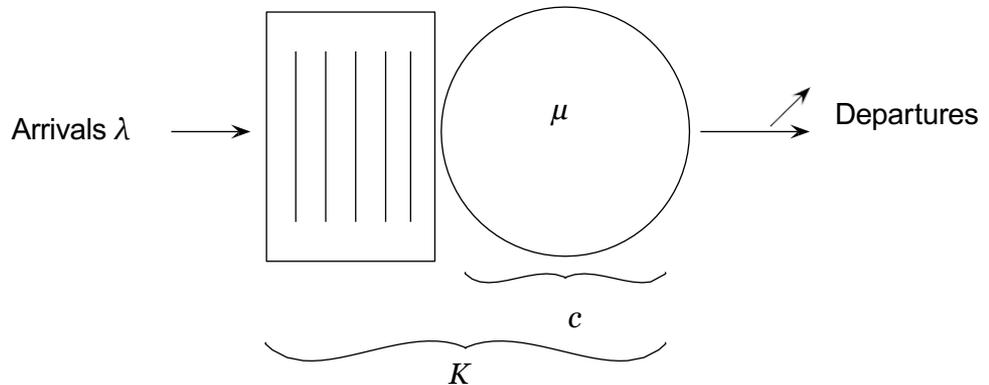
# Queueing models

- Customers requiring service are generated over time by an input source
- These customers enter the *queueing system* and join a queue
- At certain times, a member of the queue is selected for service by some rule known as the *queue discipline*.
- The required service is then performed for the customer by the *service mechanism*, after which the customer leaves the queueing system.



# Queueing models

- Parameters that characterize a queue
  - Number of parallel servers,  $c$
  - Capacity,  $K$  (equal to buffer + servers, may be infinite)
  - Arrival rate,  $\lambda$
  - Service rate of one server,  $\mu$
  - Transition probabilities,  $p_{ij}$
- Arrival distribution
- Queue discipline
- Service distribution



# Complete Kendall notation

$$A / S / c / K / P / QD$$

- $A$ : inter-arrival time distribution
- $S$ : service time distribution
- $c$ : number of servers
- $K$ : total system size ( $\infty$ )
- $P$ : population size ( $\infty$ )
- $QD$ : Queue discipline (FIFO)

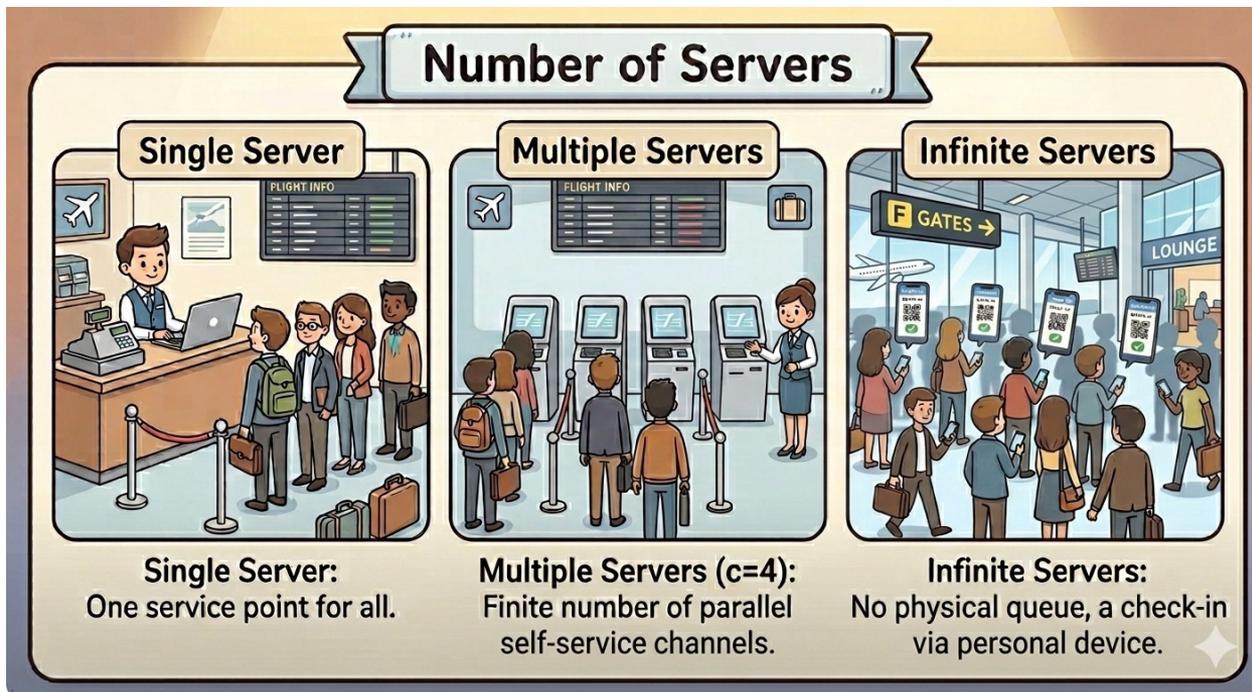
# Kendall notation

$$A / S / c / K / P / QD$$

- Arrival ( $A$ ) / Service ( $S$ ) Process
  - Assumption: i.i.d
- Some standard code letters for  $A$  and  $S$ :
  - $M$ : Exponential ( $M$  stands for memoryless/Markovian)
  - $D$ : Deterministic
  - $E_k$  :  $k$ th-order Erlang distribution
  - $G$ : General distribution
- Examples:
  - $D/D/1$ , lends itself to a graphical analysis (Unit 1)
  - $M/M/c$

# Number of servers

- **Single server:** One server for all queued customers
- **Multiple servers:** Finite number of “identical” servers operating in parallel
- **Infinite servers:** A server for every customer



# Kendall notation

$$A / S / c / K / P / QD$$

- $K$ : total system size, i.e. buffer size + number of servers
- Referred to as **capacity** in queueing theory
- $K < \infty$ : finite capacity queues

*Couldn't AI-generate an illustrative cartoon 😊 – maybe next time*

# Queue discipline

- Refers to the **order** in which members of the queue are selected for service
- **FIFO**: first-in first-out (a.k.a. FCFS); first member to arrive is first to depart, no passing
  - Single road lane, airport check-in counters
- **LIFO**: last-in first-out; last customer into queue is first to leave
  - Unboarding cars from a ferry, unboarding a bus from behind
- **Priority**: members get served in order of priority (highest to lowest)
  - Flight departures along a runway, priority seating when boarding flights
  - Yields / intersections: priority between approaches
- **SIRO**: service in random order
  - Claiming baggage at airport; hailing a taxi during rush hour
- **PS**: processor sharing; all members receive service simultaneously
  - Highway traffic flow: Road capacity is shared among vehicles
  - Airspace sector capacity: Aircraft within control sector share controller attention
- **FIFO is the most common discipline in transportation**
  - **FIFO + PS is common under congested settings**

# Outline

1. Fundamental queueing models
2. **Stationary analysis and Little's law**
3. M/M/1: Detailed analysis
4. More queues
5. Airport security checkpoints problem

# Stationary analysis

- **State of system**: number of customers in the system  $n$
- **Steady state condition**: system is independent of initial state and has reached its long-term equilibrium characteristics
  - A.k.a. steady state regime, stationary regime
- Given:
  - $\lambda$  = arrival
  - $\mu$  = service rate per server
  - $c$  = number of servers (parallel service channels)
- Quantities of interest:
  - $\bar{N}$ : expected number of users in queueing system ( $\bar{N} = E[N]$ )
  - $\bar{N}_q$ : expected number of users in queue ( $\bar{N}_q = E[N_q]$ )
  - $\bar{T}$ : expected time in queueing system per user ( $\bar{T} = E[T]$ )
  - $\bar{T}_q$ : expected waiting time in queue per user ( $\bar{T}_q = E[T_q]$ )
- 4 unknowns  $\implies$  need 4 equations
- Also of interest:  $(P_n)$ : **stationary queue length distribution**
  - $\sum_{i=0}^{\infty} P_i = 1$

# Stability

- A system is said to be **stable** if its long run averages  $(N, T)$  exist and are finite

- Consider an infinite capacity queue:

- Traffic intensity (also called utilization factor):

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

- $c\mu$ : queue service rate.
  - The queue is stable if and only if  $\rho < 1$
  - If a system is unstable, its long run measures are meaningless
  - Note:
    - This is necessary only for infinite capacity queues
    - Finite capacity queues have bounded queue lengths, and are therefore always stable
    - Stable systems  $\rightarrow$  a steady state condition exists

# Little's law

- Recall deterministic version from Unit 1
- John Little, MIT Institute Professor
- Proof in: “A proof for the queuing formula:  $L = \lambda W$ ” (1961), Operations Research
- [Little's Law as viewed on its 50th Anniversary](#) (INFORMS)
- $\bar{N} = \lambda \bar{T}$  (1)
  - $\bar{N}$ : expected number of vehicles in the system
  - $\lambda$ : system arrival rate
  - $\bar{T}$ : expected time in the system
- Assumption: The system is in a stationary regime
- No assumptions/restrictions on the:
  - inter-arrival and service time distributions
  - queue discipline
  - number of servers
- Little's law applies to several classes/categories of users
- Stationarity is not required when considering a finite time horizon (i.e., operating hours over 1 day).

# Little's law

$$\bar{N} = \lambda \bar{T} \quad (1)$$

- $\bar{N}$ : expected number of vehicles in the system
- $\lambda$ : system arrival rate
- $\bar{T}$ : expected time in the system

$$\bar{N}_q = \lambda \bar{T}_q \quad (2)$$

- $\bar{N}_q$ : expected number of vehicles in the buffer
- $\lambda$ : system arrival rate
- $\bar{T}_q$ : expected time in the buffer

# Relationships between $\bar{N}$ , $\bar{N}_q$ , $\bar{T}$ , and $\bar{T}_q$

- Little's law:
  - $\bar{N} = \lambda \bar{T}$  (1)
  - $\bar{N}_q = \lambda \bar{T}_q$  (2)
- $\bar{T} = \bar{T}_q + \frac{1}{\mu}$  (3)
  - $\mu = \text{service rate (Hz)} \Rightarrow \text{expected service time} = \frac{1}{\mu}$
- $\bar{N} - \bar{N}_q = \frac{\lambda}{c\mu}$  (for M/M/c) (4)
  - which represents the expected number of vehicles under service (in steady-state)
- Obtain one of the performance measures, the other three can then be deduced
- Let's try to obtain  $\bar{N}$ .
  - The determination of  $\bar{N}$  may be hard or easy depending on the type of queueing model at hand
  - It is easy for M/M/1 and quite easy for M/M/c and for M/G/1
- In general:  $\bar{N} = \sum_{n=0}^{\infty} nP_n$ , where  $P_n$  is the probability that there are  $n$  customers in the system

# Outline

1. Fundamental queueing models
2. Stationary analysis and Little's law
3. **M/M/1: Detailed analysis**
4. More queues
5. Airport security checkpoints problem

# Analysis of queueing models

- Closed-form expressions for the main performance measures typically involve:
  - stationary regime (i.e. steady state analysis)
  - specific distributional assumptions
- $M/M/1$  queueing system: “simple” to analyze
- General strategy:
  - Compute steady state probabilities  $P_n$
  - Compute  $\bar{N} = \sum_{n=0}^{\infty} nP_n$
  - Obtain  $\bar{N}_q$ ,  $\bar{T}$ , and  $\bar{T}_q$
- Computational techniques allow us to numerically evaluate performance measures for more general queues, and also for transient regime (i.e. dynamic analysis)

# Detailed analysis of $M/M/1$ queueing system

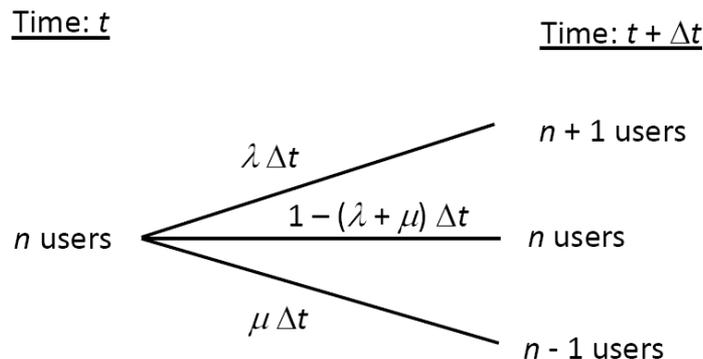
- (Recall) Inter-arrival times:

$$f_X(t) = \lambda e^{-\lambda t} \quad t \geq 0; \quad E[X] = \frac{1}{\lambda}; \quad \sigma_X^2 = \frac{1}{\lambda^2}$$

- (Recall) Service times:

$$f_S(t) = \mu e^{-\mu t} \quad t \geq 0; \quad E[S] = \frac{1}{\mu}; \quad \sigma_S^2 = \frac{1}{\mu^2}$$

- From the properties of exponential r.v.'s, the probabilities of transitions in the next  $\Delta t$  (assume small):

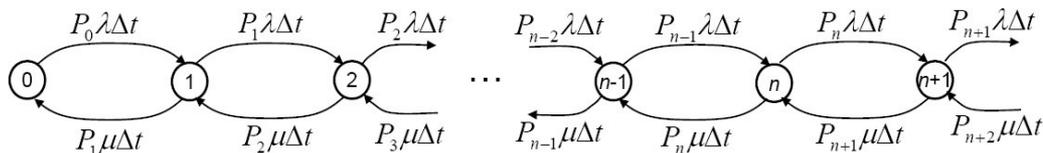


# State transition diagram for $M/M/1$

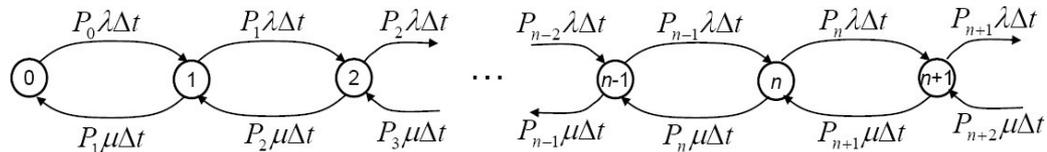
- States (number of “customers” in the system):



- The probability of observing a transition from state  $i$  to state  $j$  during the next  $\Delta t$  with the system in steady-state:

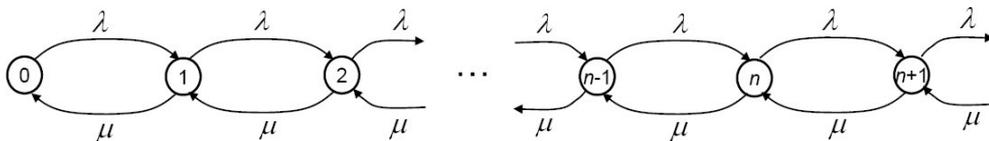


# State transition diagram for $M/M/1$



■ Another way to represent this **State transition diagram**:

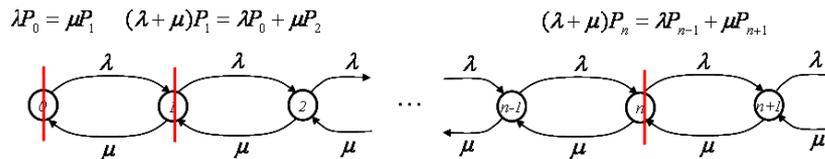
- Nodes: states
- Arcs: possible state transitions



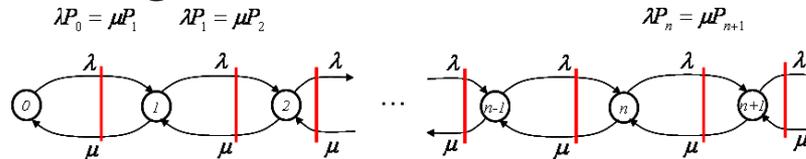
# Observing the diagram from two points

In steady state:

- At a state: flows in = flows out



- Between states: net change = 0



- The two sets of equations yield the same solutions

## $M/M/1$ : deriving $P_0$ and $P_n$

$$1. \quad P_1 = \frac{\lambda}{\mu} P_0, \quad P_2 = \left(\frac{\lambda}{\mu}\right)^2 P_0, \dots, \quad P_n = \left(\frac{\lambda}{\mu}\right)^n P_0$$

$$2. \quad \sum_{n=0}^{\infty} P_n = 1, \Rightarrow P_0 \sum_{n=0}^{\infty} \left(\frac{\lambda}{\mu}\right)^n = 1, \Rightarrow P_0 = \frac{1}{\sum_{n=0}^{\infty} \left(\frac{\lambda}{\mu}\right)^n}$$

$$3. \quad \text{For } |x| < 1 \text{ (stable queues), } \sum_{n=0}^{\infty} x^n = \frac{1}{1-x}$$

$$4. \quad \text{Define: } \rho = \frac{\lambda}{\mu}$$

$$P_0 = \frac{1}{\sum_{n=0}^{\infty} \rho^n} = 1 - \rho$$

$$P_n = \rho^n (1 - \rho)$$

# $M/M/1$ : deriving $\bar{N}$ , $\bar{N}_q$ , $\bar{T}$ , and $\bar{T}_q$

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

$$\bar{T} = \frac{\bar{N}}{\lambda} = \frac{\lambda}{\mu-\lambda} \cdot \frac{1}{\lambda} = \frac{1}{\mu-\lambda}$$

$$\bar{T}_q = \bar{T} - \frac{1}{\mu} = \frac{1}{\mu-\lambda} - \frac{1}{\mu} = \frac{\lambda}{\mu(\mu-\lambda)}$$

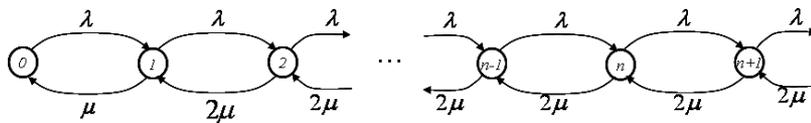
$$\bar{N}_q = \lambda \bar{T}_q = \frac{\lambda^2}{\mu(\mu-\lambda)}$$

# Outline

1. Fundamental queueing models
2. Stationary analysis and Little's law
3. M/M/1: Detailed analysis
4. **More queues**
  - a.  $M/M/2$
  - b.  $M/M/c$
  - c.  $M/M/c/K$
5. Airport security checkpoints problem

# M/M/2 queueing system

- What happens if we have two parallel, independent servers, each with service rate  $\mu$  and exponentially distributed service times?
- State transition diagram:



- Balance equations: 
$$\begin{cases} \lambda P_0 = \mu P_1 \\ \lambda P_1 = 2\mu P_2 \\ \lambda P_2 = 2\mu P_3, \dots \end{cases}$$

$$\begin{cases} P_1 = \frac{\lambda}{\mu} P_0 \\ P_2 = \frac{\lambda}{2\mu} P_1 = \frac{1}{2} \left(\frac{\lambda}{\mu}\right)^2 P_0 \\ P_3 = \frac{\lambda}{2\mu} P_2 = \left(\frac{1}{2}\right)^2 \left(\frac{\lambda}{\mu}\right)^3 P_0 \\ P_n = \left(\frac{1}{2}\right)^{n-1} \left(\frac{\lambda}{\mu}\right)^n P_0 \end{cases}$$

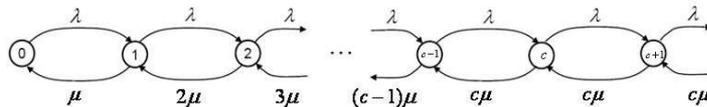
# M/M/2 queue

$$P_0 = \frac{1 - \frac{\lambda}{2\mu}}{1 + \frac{\lambda}{2\mu}} \quad \text{for } \lambda < 2\mu$$

$$P_n = 2 \left( \frac{\lambda}{2\mu} \right)^n P_0, \quad n \geq 1$$

$$\dots \bar{N} = \frac{\frac{\lambda}{\mu}}{\left(1 + \frac{\lambda}{2\mu}\right) \left(1 - \frac{\lambda}{2\mu}\right)}$$

# M/M/c queue



- Most common queuing model
- Example: airport check-in counters
- Traffic intensity / utilization factor:  $\rho = \frac{\lambda}{c\mu}$
- Stability:  $\frac{\lambda}{c\mu} < 1$
- Stationary dbn:

$$P_k = \begin{cases} \frac{(\lambda/\mu)^k}{k!} P_0, & k = 1, 2, \dots, c-1 \\ \frac{(\lambda/\mu)^k}{c! c^{k-c}} P_0, & k = c, c+1, \dots \end{cases}$$

$$P_0 = \left[ \frac{(\lambda/\mu)^c}{c! (1-\rho)} + \sum_{k=0}^{c-1} \frac{(\lambda/\mu)^k}{k!} \right]^{-1}$$

- Expected queue length (in the buffer)?

# $M/M/c$ queue

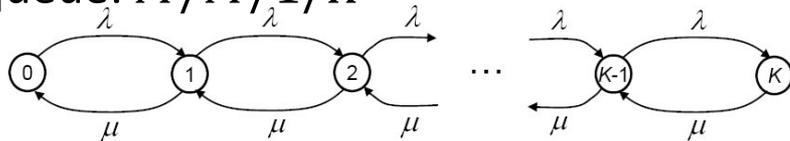
- Probability that an arrival must wait
  - iff upon arrival there are  $\geq c$  customers
  - $P(\text{wait}) = \frac{(\lambda/\mu)^c}{c!(1-\rho)} P_0$  [Erlang-C formula]

Can derive the usual quantities:

- Little's formula:  $T_q = \frac{N_q}{\lambda}$
- $T = T_q + \frac{1}{\mu}$
- To obtain  $N$ :
  1. Little's formula:  $N = \lambda T$
  2.  $N = N_q + \frac{\lambda}{\mu}$

# $M/M/c/K$

- Finite capacity queue:  $M/M/1/K$



$$P_n = \frac{1 - \rho}{1 - \rho^{K+1}} \rho^n \quad n = 0, 1, \dots, K$$

- The queue is always stable ( $\forall \rho$ ), so steady state is always reached
- When system is full: arrivals are lost
- Effective* arrival rate (i.e. arrivals that actually enter system):  $\lambda(1 - P_K)$
- Careful when applying Little's Law!** Count only the vehicles that actually join the system:

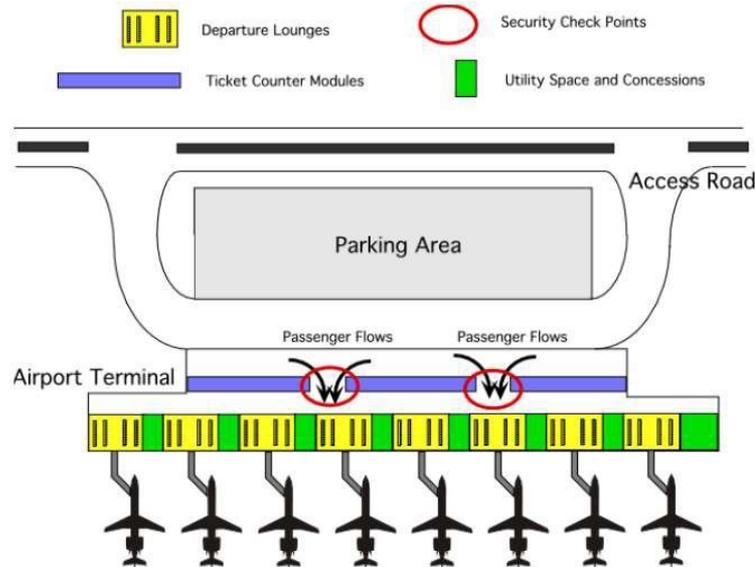
$$\begin{aligned} \lambda' &= \lambda(1 - P_K) \\ \bar{N} &= \lambda' \bar{T} \end{aligned}$$

# Outline

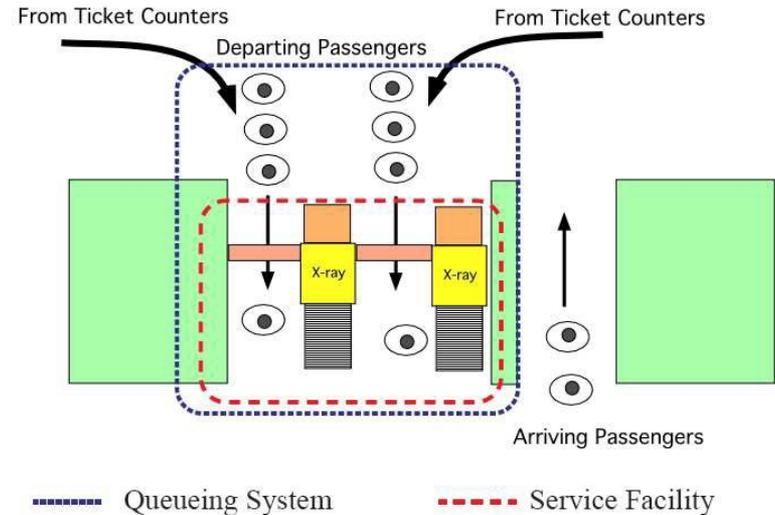
1. Fundamental queueing models
2. Stationary analysis and Little's law
3. M/M/1: Detailed analysis
4. More queues
5. **Airport security checkpoints problem**

# Example: Airport security checkpoints

- Level of service at airport security checkpoints
- Airport terminal layout
- An airport has two security checkpoints for all passengers boarding aircrafts.



- Security checkpoint layout
- Each security check point has two x-ray machines.



# Airport security checkpoints

## Setting

- A survey reveals that on average a passenger takes 45 seconds to go through the x-ray machine.
- The service times are exponentially distributed.
- The arrivals to a checkpoint follow a Poisson distribution with a constant arrival rate of one passenger every 25 seconds.
- The demand for services in 2035 is expected to grow by 60% compared to that of today.

## Problem

1. What is the current utilization of a checkpoint (i.e., two x-ray machines)?
2. What should be the number of x-ray machines for the design year of this terminal (year 2035), if the expected waiting time is to be kept under 2 minutes?
3. What is the new utilization of the future facility?
4. What is the probability that more than 4 passengers wait for service in the design year?

# Airport security checkpoints

## Setting

- A survey reveals that on average a passenger takes 45 seconds to go through the x-ray machine.
- The service times are exponentially distributed.
- The arrivals to a checkpoint follow a Poisson distribution with a constant arrival rate of one passenger every 25 seconds.
- The demand for services in 2035 is expected to grow by 60% compared to that of today.

## Problem

1. What is the current utilization of a checkpoint (i.e., two x-ray machines)?
2. What should be the number of x-ray machines for the design year of this terminal (year 2035), if the expected waiting time is to be kept under 2 minutes?
3. What is the new utilization of the future facility?
4. What is the probability that more than 4 passengers wait for service in the design year?

# Airport security checkpoints

## Setting

- A survey reveals that on average a passenger takes 45 seconds to go through the x-ray machine.
- The service times are exponentially distributed.
- The arrivals to a checkpoint follow a Poisson distribution with a constant arrival rate of one passenger every 25 seconds.
- The demand for services in 2035 is expected to grow by 60% compared to that of today.

## Problem

1. What is the current utilization of a checkpoint (i.e., two x-ray machines)?
2. What should be the number of x-ray machines for the design year of this terminal (year 2035), if the expected waiting time is to be kept under 2 minutes?
3. What is the new utilization of the future facility?
4. What is the probability that more than 4 passengers wait for service in the design year?

# References

1. Larson, Richard C. and Amedeo R. Odoni. Urban Operations Research. Prentice-Hall (1981). Chapter 4: Queueing Theory.
2. John Little, Little's Law as Viewed on its 50th Anniversary. Operations Research, vol 59. 2011.  
<https://www.informs.org/Blogs/Operations-Research-Forum/Little-s-Law-as-Viewed-on-its-50th-Anniversary>
3. Slides adapted from Carolina Osorio

## *Appendix: More queues*

More tools in the stochastic modeling toolbox!

1.  $M/M/c/c$  - Erlang loss model
2.  $M/D/1$

# $M/M/c/c$ - Erlang loss model

- $M/M/c/c$ : Erlang loss model
- First queueing model to be investigated
- Agner Krarup Erlang
  - Danish telephone engineer, who investigated it in the early 1900's as a model for telephone switch which can handle only  $c$  calls
  - Queueing theory pioneer
  - The theory of stochastic processes was not yet developed at the time
  - Erlang derived a formula for the proportion of lost calls, Erlang loss formula:

$$P_c = \frac{(\lambda/\mu)^c / c!}{\sum_{i=0}^c (\lambda/\mu)^i / i!}$$



## $M/M/c/c$ - Insensitivity

- Erlang assumed exponential service times, but conjectured that it would hold for generally distributed service times
- Insensitivity: the loss probability is insensitive to the form of the service time distribution; it depends only on its expectation.
- This was not proved until the 1960's

$$P_c = \frac{(\lambda/\mu)^c / c!}{\sum_{i=0}^c (\lambda/\mu)^i / i!}$$

- Loss probability holds for  $M/G/c/c$  queues

## $M/G/c/c$ - Erlang loss model

- Actually, the insensitivity property also holds for the stationary distribution
- Insensitivity: the stationary distribution is insensitive to the form of the service time distribution; it depends only on its expectation.

$$P_n = \frac{(\lambda/\mu)^n/n!}{\sum_{i=0}^c (\lambda/\mu)^i/i!}, \forall n \in [0, c]$$

- Insensitivity property  $\rightarrow$  Erlang loss model is of wide interest
  - Model is commonly used for the analysis of telecommunication systems, also: urban service systems, inventory, reliability

# $M/D/1$ queue

- Has been used to model vehicles on a lane at signalized urban intersections
- Exponentially distributed inter-arrival times
- Deterministic service distribution
- One server
  
- Recall the traffic intensity:  $\rho = \frac{\lambda}{\mu}$ 
  - $\rho$ : traffic intensity
  - $\lambda$ : arrival rate [veh/unit time]
  - $\mu$ : service rate [veh/unit time]

# $M/D/1$ queue

- For a stable queue ( $\rho < 1$ ):
  - Expected number of vehicles in the buffer [veh]:

$$N_q = \frac{\rho^2}{2(1 - \rho)}$$

- Expected waiting time in the buffer (per veh)

$$T_q = \frac{\rho}{2\mu(1 - \rho)}$$

- Expected time in the system: sum of the expected waiting time and the expected service time:

$$T = \frac{2 - \rho}{2\mu(1 - \rho)}$$

- **Note:** traffic intensity:  $\rho < 1$ , then:
  - the  $D/D/1$  queue predicts no queue formation,
  - models with probabilistic arrivals/departures (e.g.  $M/D/1$ ) predict queue formations under such conditions.

# Queueing theory - keep in mind

- Queueing theory can provide **insights** and approximation of the main system performance measures.
  - Can enable identification of the location of bottlenecks in networks,
  - Give indications on how to improve the system's performance.
- Most closed-form results involve **stationary regime** (steady-state) and **low-order moments** (mean, variance) of the inter-arrival and service time distributions
- Trade-off: realistic model (few available results) vs. tractability (assumptions are questionable)