

# Client-Server Systems

## Performance Modeling

Engineering Cloud Computing

Week 4

## Client-Server Architecture

- **Clients:** Remote programs and users
- Servers can support multiple clients
- Servers implement one or more services
- Centralized services implemented on single servers
- Distributed: service implemented across multiple servers

# Server Capacity

- Servers can become a bottleneck with increasing number of requests
  - Computational capacity, limited by the CPUs
  - Storage capacity: I/O transfer rate
  - Network between user and server

# Computation Inside Servers

- Packet processing
  - Data from NIC to CPU via DMA
  - Interrupt handling
  - Packet travels up the network stack
  - Processes blocked on socket woken up
- Application level processing
  - Parse data from socket
  - Process/store data

# Server performance

## Key Metrics

**Response time:** Time between request initiation and response

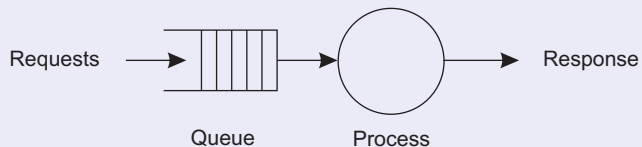
**Throughput:** Number of requests handled by server (per second)

## Server Performance Considerations

- How many concurrent clients can be served?
- What is the maximum throughput that can be sustained?
- What is the response times for clients?

# Queuing Theory Model

A centralized service can be modeled as a simple queuing system



## Assumptions and notations

- The queue has infinite capacity  $\Rightarrow$  arrival rate of requests is not influenced by current queue length or what is being processed.
- Arrival rate of requests:  $\lambda$
- Processing capacity service:  $\mu$  requests per second

## Quick Quiz

Bob finds his friend, Alice, at the bus-stop. It turns out both Alice and Bob are waiting for the same bus.

Alice has been waiting for the bus for 10 minutes. The bus is scheduled to arrive every 30 minutes.

Assume that there is no other information available about the bus (no real-time GPS etc.). What is Bob's expected waiting time for the bus?

## Distribution of Requests and Service Times

- Requests arrive according to a random process
- Typically, arrival process is modeled as a **Poisson distribution**
- Arrival rate:  $\lambda$  per second
- Request service rate:  $\mu$  per second

$$P(n \text{ arrivals in interval } T) = \frac{(\lambda T)^n e^{-\lambda T}}{n!}$$

$$E[n] = \lambda T$$

**Inter-arrival time:** Time between successive events

$$P(IA \leq t) = 1 - P(IA > t) \quad (1)$$

$$= 1 - P(0 \text{ arrivals in time } t) \quad (2)$$

$$= 1 - e^{-\lambda t} \quad (3)$$

(4)

$1 - e^{-\lambda t}$  is the CDF of the **exponential distribution!**

**Service Time:** Exponentially distributed with parameter  $\mu$



# Exponential Distribution

- CDF:  $F(t) = 1 - e^{-\lambda t}$
- Probability distribution:  $f(t) = \lambda e^{-\lambda t}$
- **Memoryless**:  $P(X \leq T + a | X > a) = P(X \leq T)$

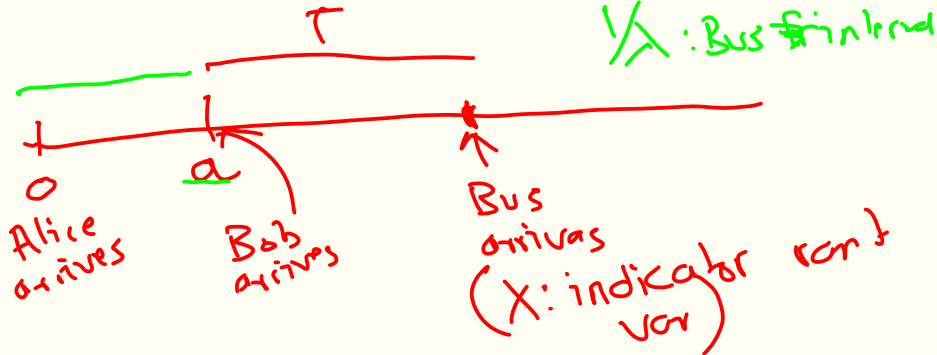
Proof:

$$P(X \leq T + a | X > a) = \frac{P(a \leq X \leq T + a)}{P(X > a)} \quad (5)$$

$$= \frac{\int_a^{T+a} \lambda e^{-\lambda t} dt}{\int_a^{\infty} \lambda e^{-\lambda t} dt} \quad (6)$$

$$= 1 - e^{-\lambda T} = P(X \leq T) \quad (7)$$

- Previous history does not help in predicting future events
- Waiting for the bus example: Waiting time of others doesn't matter if bus arrivals are exponentially distributed.



# M/M/1 Queue Properties

- $\rho = \lambda/\mu$
- $\lambda < \mu$
- Probability that system is idle =  $1 - \rho$
- Utilization =  $\rho$
- Mean number of objects in the system =  $\rho/(1 - \rho)$

Fraction of time having  $k$  requests in the system

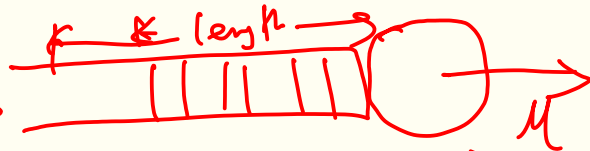
$$p_k = \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^k$$

$\uparrow$   
 $p(\text{idle})^k = \rho^k$

if  $\lambda > \mu$ :

$$1 - \frac{\lambda}{\mu} = P(\text{idle})$$

$\lambda \approx \mu$



Report interval time

$$\text{Utiliz} = 1 - P(\text{idle})$$

$$\frac{\frac{\lambda}{\mu}}{1 - \frac{\lambda}{\mu}}$$

# Little's Law ← General

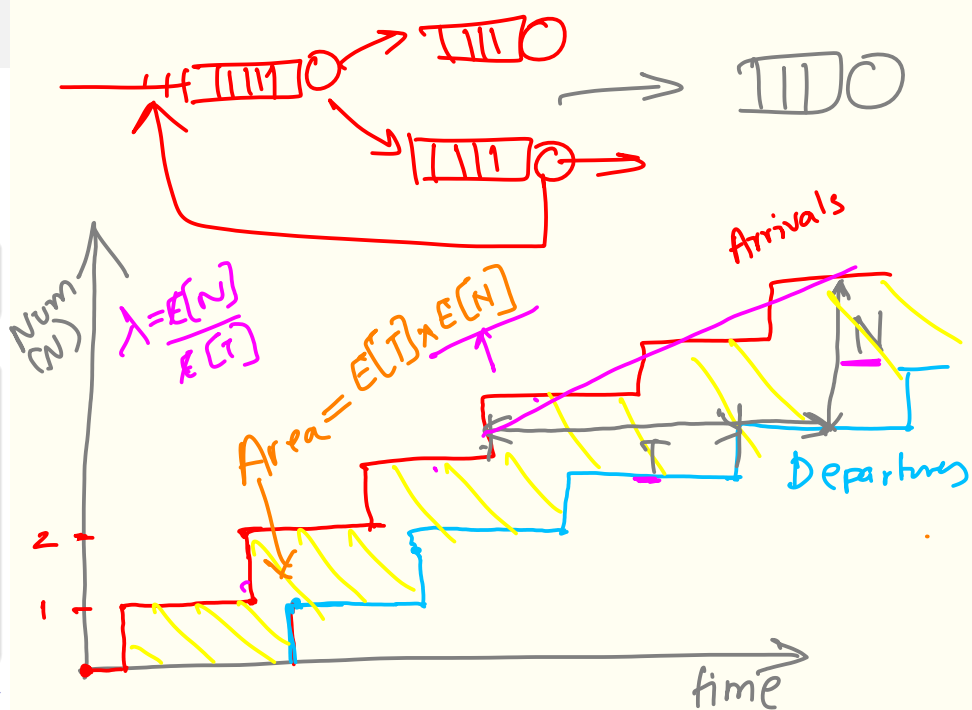
- N: Number of items in the system
- S: Response time (Time to leave the system)  $\equiv T$

## Little's Law

$$E[N] = \lambda E[T]$$

## Proof outline

- 1 Plot N vs. time, for a total time T
- 2 Area of the 'ribbon', A = Time spent by all items
- 3  $\lambda = N/T$
- 4 Mean time spent in the system,  $E[T] = A/N$
- 5 Mean number in the system,  $E[N] = A/T$
- 6 Counting the area in two ways =  $E[N] = \lambda E[T]$



## Little's Law Applications

- Very general. Multiple queueing disciplines, network of queues, etc!
- Average response time of server,  $E[S] = E[T] = E[N]/\lambda$
- $E[N] = \sum_0^{\infty} k P_k = \rho / 1 - \rho$
- $E[S] = \frac{1}{\mu - \lambda}$

### Scenarios

- Is it better to have one queue with 2 servers or 2 separate queues?
- What happens when processing power is doubled?

$P_k$ : Prob of  $k$  items

$$\sum k(1-\rho) \cdot \rho^k = (1-\rho) \sum_0^{\infty} k \rho^k$$

Little's Law  $E[S] = \frac{E[N]}{\lambda} = \frac{\rho / \mu}{1 - \rho} = \frac{\rho}{\lambda(1-\rho)}$

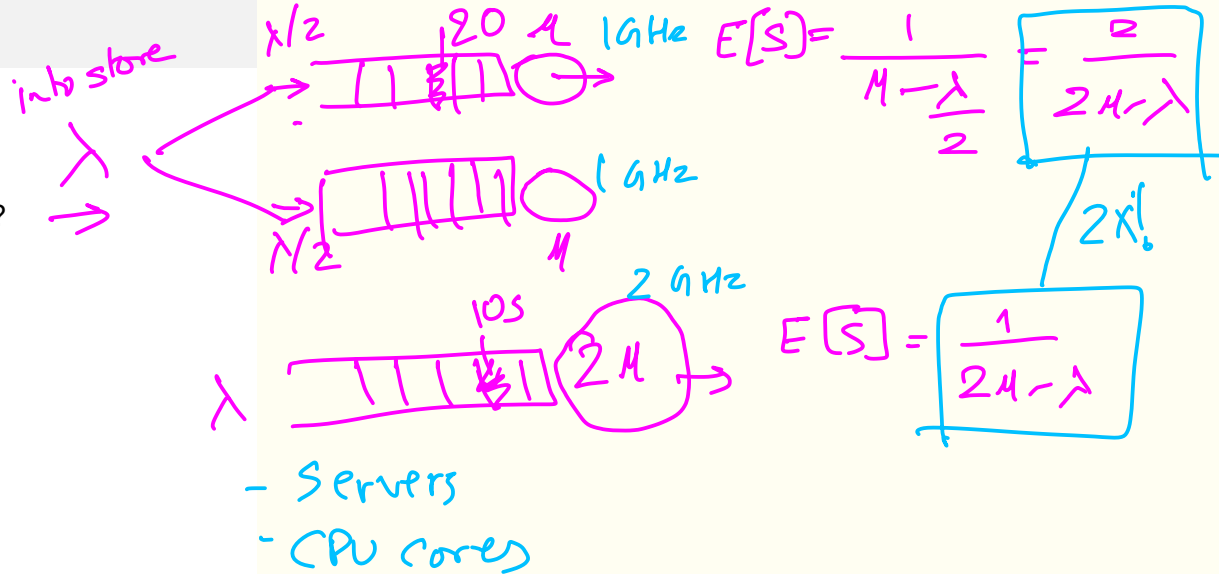
$$= \frac{\lambda}{\mu - \lambda} = \frac{1}{\mu - \lambda} = E[S]$$

# Server Performance Implications

- Average Response Time =  $1/\mu - \lambda$
- Useful to identify saturating load
- What to do if incoming traffic rate ( $\lambda$ ) is close to  $\mu$ ?
  - Scaling techniques. Next class!

$$\lambda < \mu$$

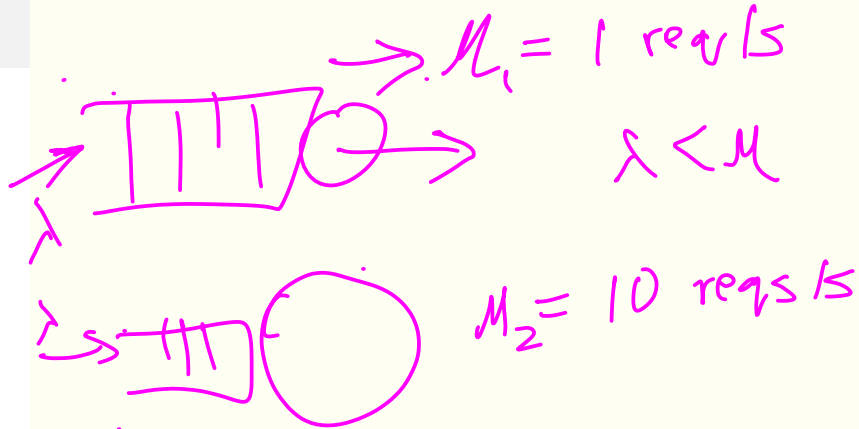
$$E[S] = \frac{1}{\mu - \lambda}$$



# Throughput

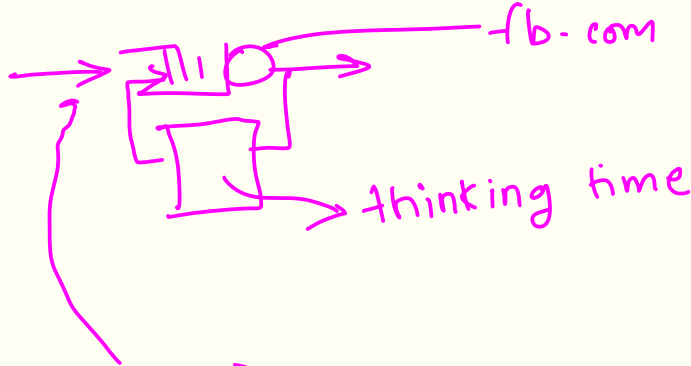
- $X$ : Rate at which events are processed
- $C$  events processed in total time  $T$
- $X = \frac{C}{T}$
- Events are only processed if system is busy
- Rate at which events are processed when system is busy =  $\mu$
- $X = \rho \cdot \mu = \lambda$
- Independent of  $\mu$ !
- Throughput of server doesn't improve if its performance improves!?!  
*Why do we want faster servers?*

$$= 1 - P(\text{not idle}) \\ = \rho = \frac{\lambda}{\mu}$$



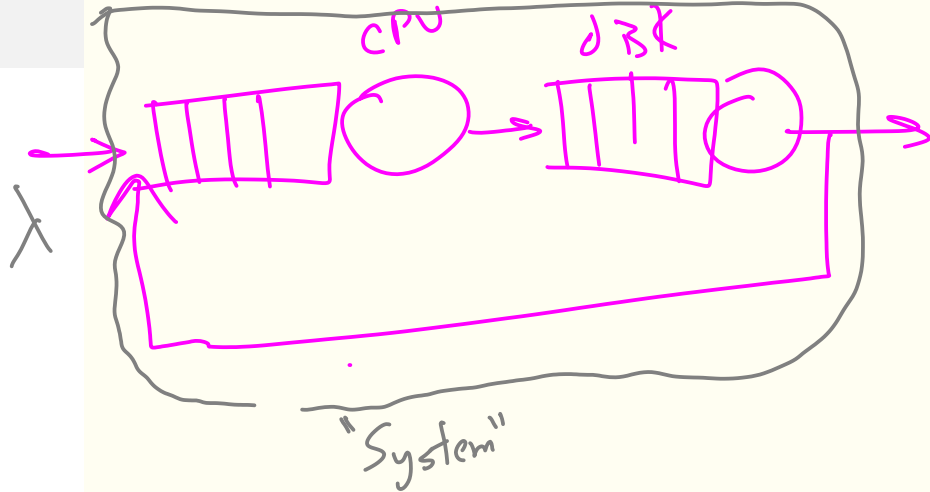
## Closed-loop Systems

- Looked at **open** systems so far
- Many systems are closed, or at least have some feedback
- Processed items feed back into the queue
- Web server example: people view a sequence of web-pages, based on what is served



## Queue Networks

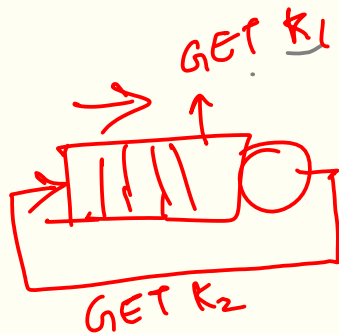
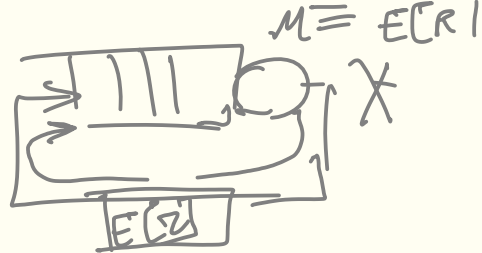
- Can represent system as a network of queues
- One queue for CPU, one for disk, etc.
- Or for different parts of the application
- Little's law is applicable!





## Closed Networks

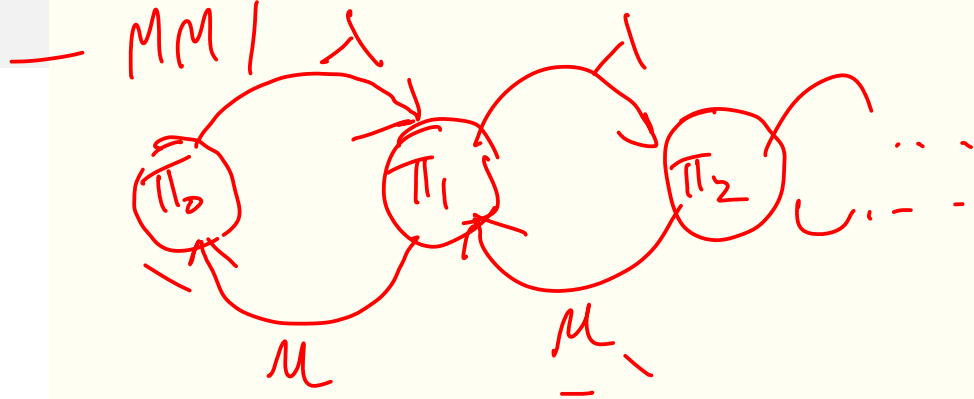
- Items feed back into queue after some “thinking time”  $E[Z]$
- Total number of items =  $N$
- $E[R]$  is the response time
- Little's law:  $N = XE[T]$
- But  $E[T] = E[R] + E[Z]$
- Throughput  $X \leq \frac{N}{E[R] + E[Z]}$
- For small  $N$ , the equality holds
- In practice, throughput converges to  $1/E[R]$  for high  $N$ . (system is saturated)
- Closed systems useful for measuring the service rate  $\mu$



while True:  
generate  
KU\_requests();

# Markov Chains

- States with transition probabilities  $P_{ij}$  between state  $i$  and  $j$
- Represented by a matrix  $\mathbf{P}$
- $P_{ij}$  is probability of going from  $j$  to  $i$  in 1 step
- $(\mathbf{P}^2)_{ij}$  denotes probability of being in  $j$  after starting at  $i$  after 2 steps.
- We are interested in  $P^n$  for  $n \rightarrow \infty$
- For markov chains,  $P_{ij}^n$  is in a row and column
- That is, the limiting probability of being in a state doesn't depend on where you start.
- Limiting distribution of being in state  $j$ :  $\pi_j = \lim_{n \rightarrow \infty} P_{ij}^n$
- $\sum \pi_i = 1$



# M/M/1 Queue Markov Chains

■ Balance equations:

↑ ■  $\lambda\pi_0 = \mu\pi_1$

■  $(\lambda + \mu)\pi_1 = \lambda\pi_0 + \mu\pi_2$

■  $(\lambda + \mu)\pi_n = \lambda\pi_{n-1} + \mu\pi_{n+1}$

$$\pi_1 = \frac{\lambda}{\mu}\pi_0$$

$$\pi_n = \left(\frac{\lambda}{\mu}\right)^n \pi_0$$

Let  $\rho = \lambda/\mu$  and  $\rho < 1$

$$1 = \sum_{i=0}^{\infty} \pi_i = \pi_0 \sum_{i=0}^{\infty} \rho^i = \frac{\pi_0}{1 - \rho} \quad (8)$$

Last step is using the geometric series sum for  $\rho$

$$\pi_0 = 1 - \rho$$

$$\text{Def: } \pi_0 \lambda$$
$$I_n = \pi_1 \mu$$

$$\pi_1 = \rho \pi_0$$
$$\pi_n = \rho^n \pi_0$$