

Operating systems

Scheduling 2: queueing theory

Last time...

- CPU scheduler
- simple scheduling policies
 - FCFS
 - SJF
 - RR
- more realistic
 - multi-level, multi-level with feedback
 - Linux scheduler

Today...

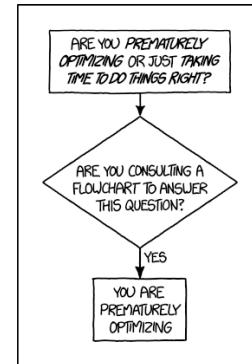
- scheduling criteria
- how to evaluate algorithms
- **queueing theory**

Scheduling criteria

- waiting time
- response time (or delay)
- throughput
- predictability
- scheduling overhead
- starvation

How to evaluate a scheduling algorithm?

1. deterministic modeling
2. queueing models
3. simulations
4. implementation



<https://xkcd.com/1691/>



queueing theory

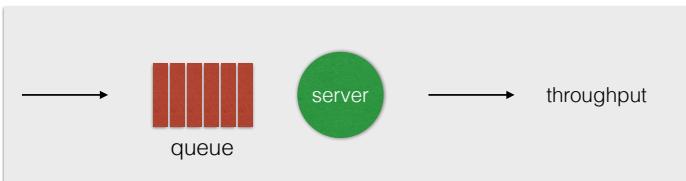
friscokids.net

We will assume

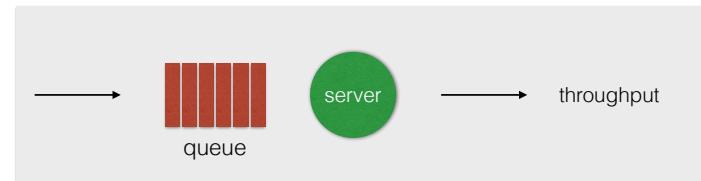
- FIFO scheduling
- all tasks that arrive will eventually be serviced



our system
is assumed to be
work-conserving



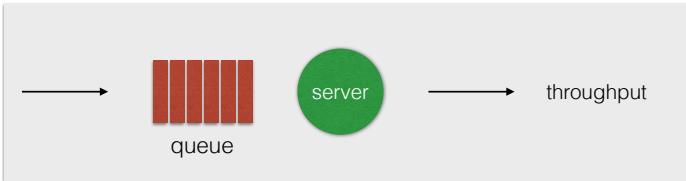
$$\text{queueing delay} + \text{service time} = \text{response time}$$



$$U = \min \left(1, \frac{\lambda}{\mu} \right)$$

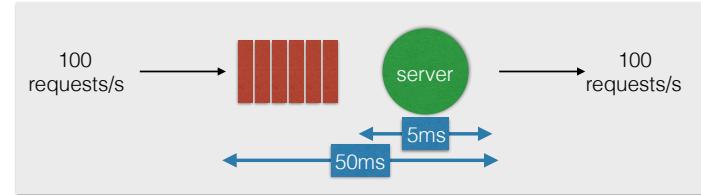
utilization

arrival rate
 $\text{service rate} = \frac{1}{\text{service time}}$



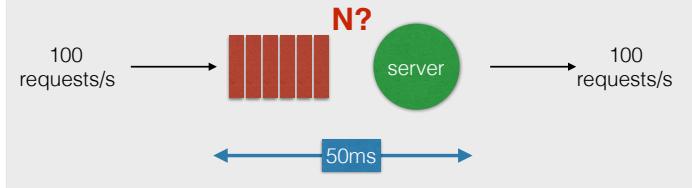
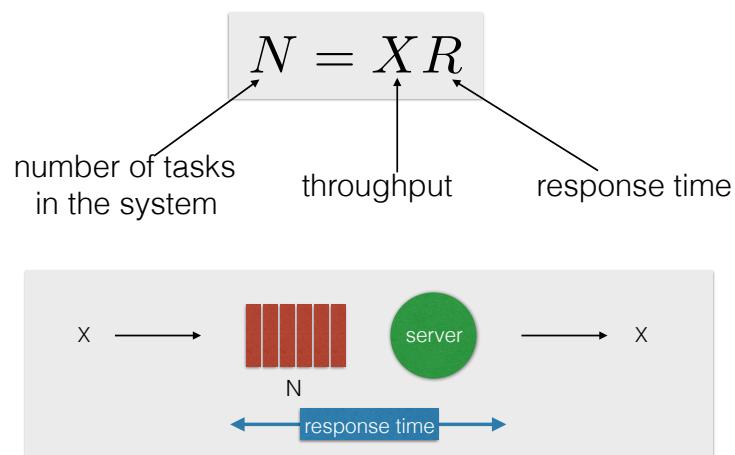
$$U = \min \left(1, \frac{\lambda}{\mu} \right)$$

$$X = U\mu$$

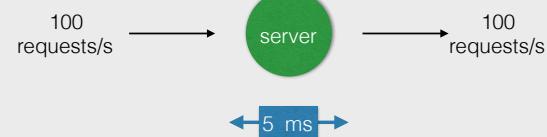


What is the queueing delay?

Little's law

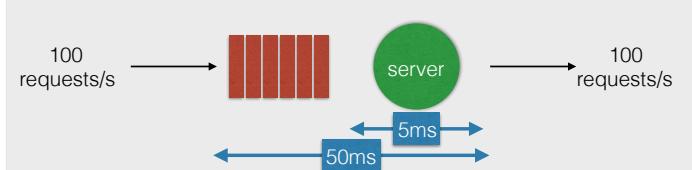


$$N = X R$$



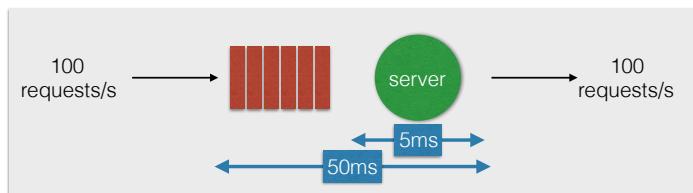
What is the average utilization of the server?

$$N = X R$$



1. How long does an average request spend in the queue?
2. How many requests are in the queue?

$$N = X R$$



Average number of tasks in queue: 4.5

execution time: 5 ms

Why is queueing delay 4.5 ms and not $4.5 \times 5 \text{ ms} = 22.5 \text{ ms}$?

$$N = X R$$

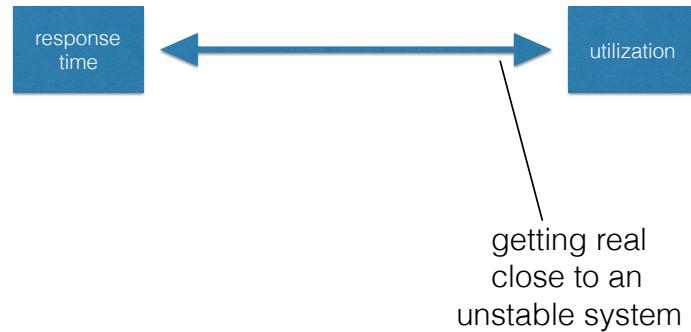


Scheduling criteria

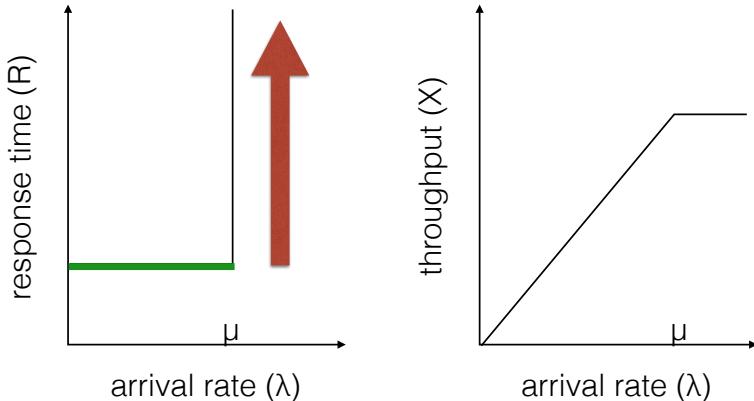
- waiting time
- response time (or delay)
- throughput
- predictability
- scheduling overhead
- starvation

what about utilization?

no free lunch

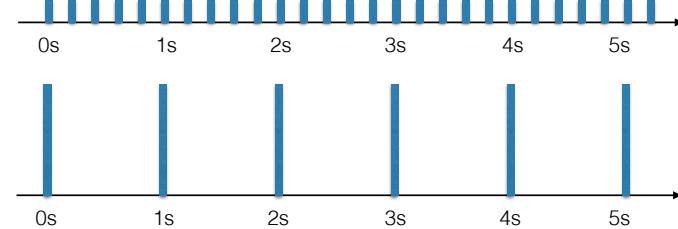


extreme scenario 1: evenly spaced arrivals

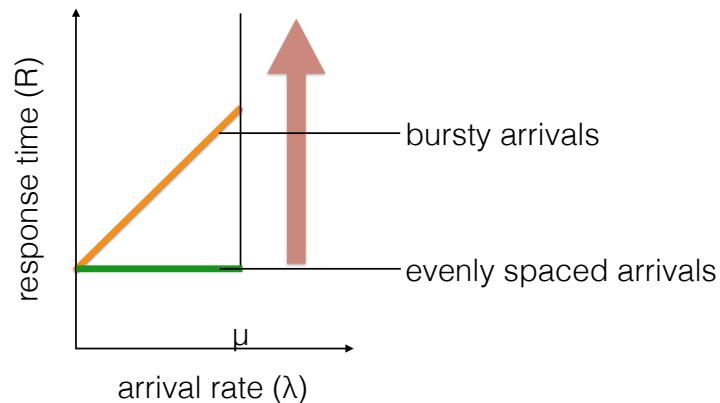


example: $\lambda = 5$ requests/s

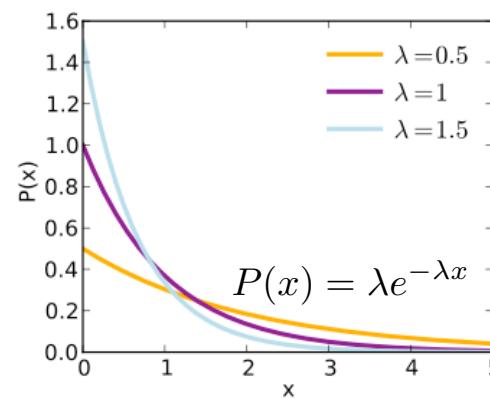
evenly spaced



extreme scenario 2: bursty arrivals



Exponential distribution

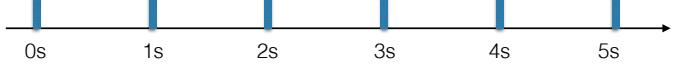


example: $\lambda = 5$ requests/s

evenly spaced



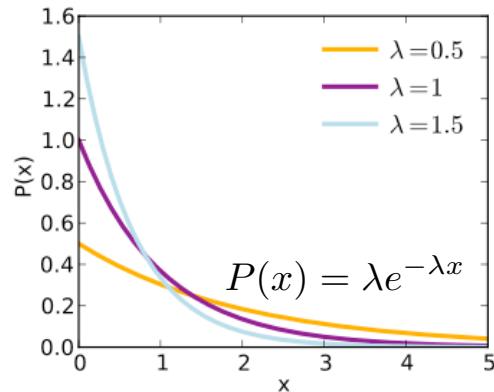
bursty



exponentially distributed



Exponential distribution



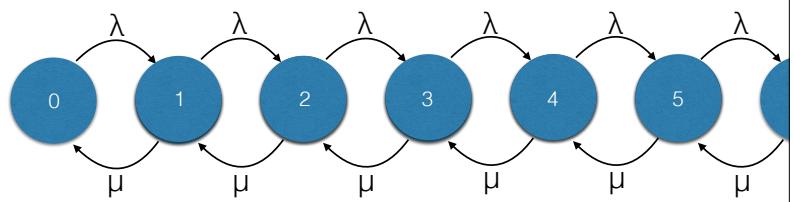
- memoryless
- light-tailed

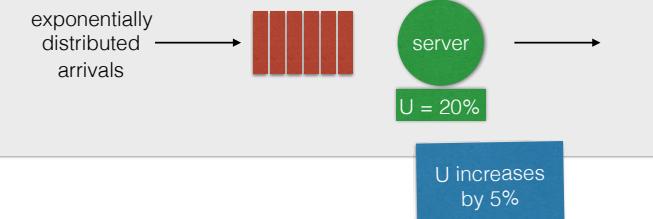
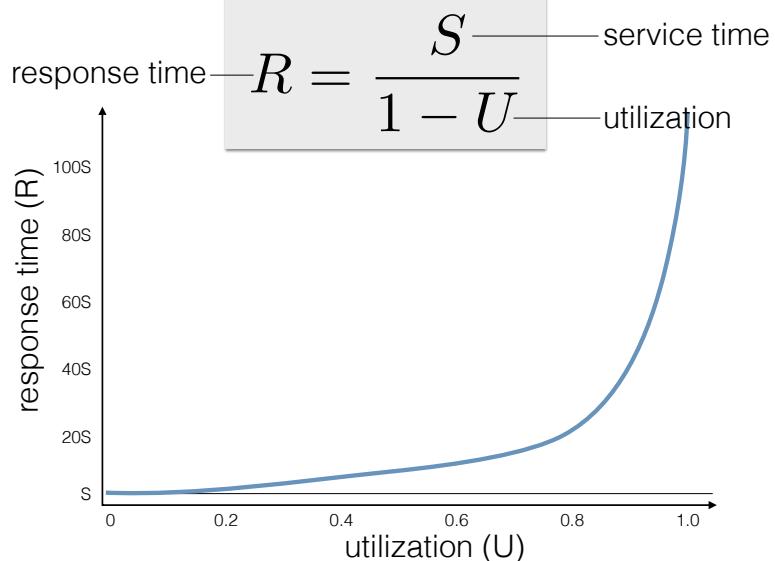


arrival time:
memory-full or
memoryless
probability distribution?



if the arrival process is a
memoryless distribution

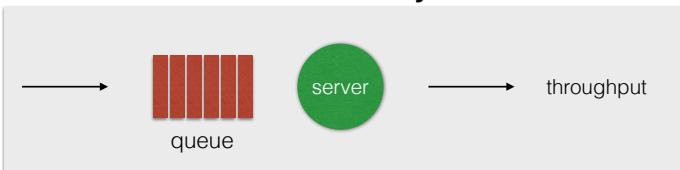




1. by how much does the response time increase?
2. how does that increase compare to the case when utilization goes from 90% to 95%?

$$R = \frac{S}{1 - U}$$

in summary...



$$\text{Little's law } N = XR$$

response time \longleftrightarrow utilization

