

CSC 2232: Topics in Computer System Performance and Reliability

Bianca Schroeder

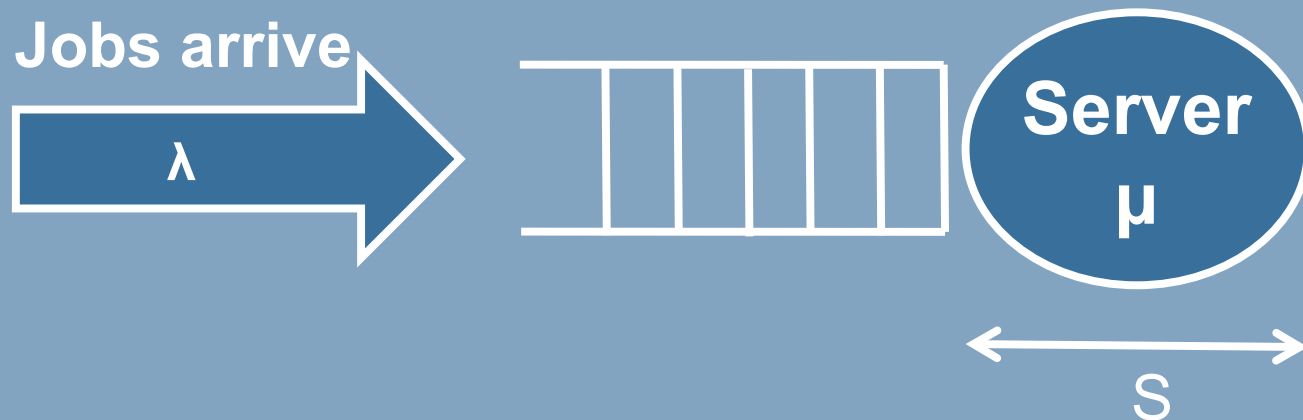
**Department of Computer Science
University of Toronto**

TODAY'S AGENDA

- More operational laws
- Modification analysis

- After this lecture, you can don a suit and call yourself a systems consultant!

THE SINGLE SERVER NETWORK



- Average arrival rate: $\lambda = 1 \text{ job} / 6 \text{ sec}$
- Mean interarrival time: $1/\lambda = 6 \text{ sec}$
- Service Requirement S
- Mean service time: $E[S] = 3 \text{ sec}$
- Average service rate: $\mu = 1/E[S] = 1 \text{ job} / 3 \text{ sec}$
- Service Order: FCFS (First-Come-First-Serve)

PERFORMANCE METRICS

- Response time: T_S
 - Turnaround time, flow time, time in system
 - Main interest: $E[T]$, also $\text{Var}[T]$
- Waiting time: T_Q
 - Time in queue
- Number of jobs in system: N_S
- Number of jobs in queue: N_Q

PERFORMANCE METRICS

Metrics:

T_S

N_S

T_Q

N_Q

- Definitions:
- U_i : Fraction of time device i is busy
- X_i : The rate of completions at device i (in jobs/ sec)
- How does X_i relate to U_i ?

HOW DOES X_i RELATE TO U_i

X_i = Mean rate of completion

= $E\{\text{Rate of completion} \mid \text{server busy}\}$
 $P\{\text{server busy}\}$

+ $E\{\text{Rate of completion} \mid \text{server idle}\}$
 $P\{\text{server is idle}\}$

$$= \mu * U_i$$

How does this change with changes in

- Job size distribution
- Interarrival time distribution
- Service order

$$X_i = \mu * U_i$$

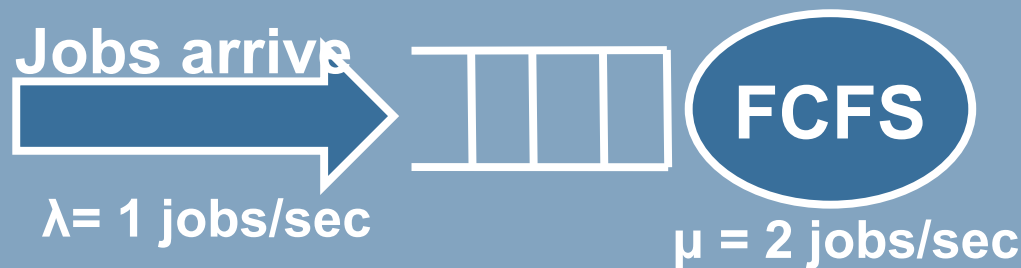
THE UTILIZATION LAW

$$X_i = \mu * U_i$$

or equivalently

$$U_i = X_i * E[S]$$

WHAT IS THROUGHPUT?



$$X_i = \mu * U_i$$

- What is U?
 - Fraction of time server is busy
= mean service time / mean time btw. Arrivals
= $(1/\mu)/(1/\lambda) = \lambda/\mu$

Throughput
does not depend
on service rate!

$$X_i = \lambda$$

BACK TO OUR OLD EXAMPLE

Metrics:

T_S
N_S
T_Q
N_Q
X
U

$$X_i = \mu * U_i = \lambda$$

Jobs arrive

 $\lambda = 1$ jobs/sec



FCFS

$\mu = 2$ jobs/sec

Jobs arrive

 $\lambda = 1$ jobs/sec

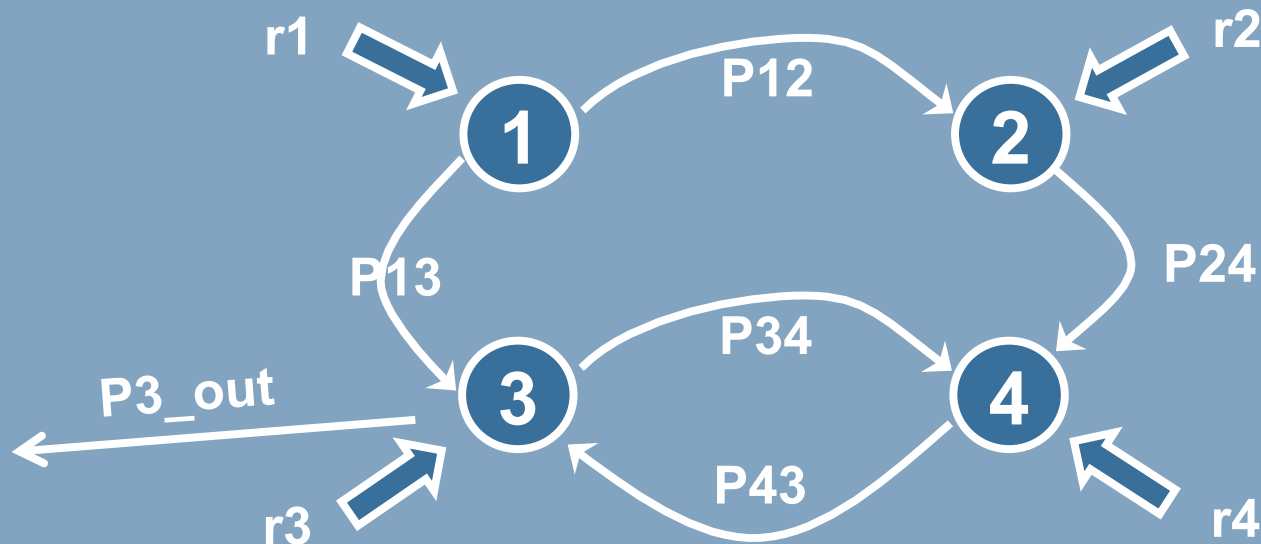


FCFS

$\mu = 4$ jobs/sec

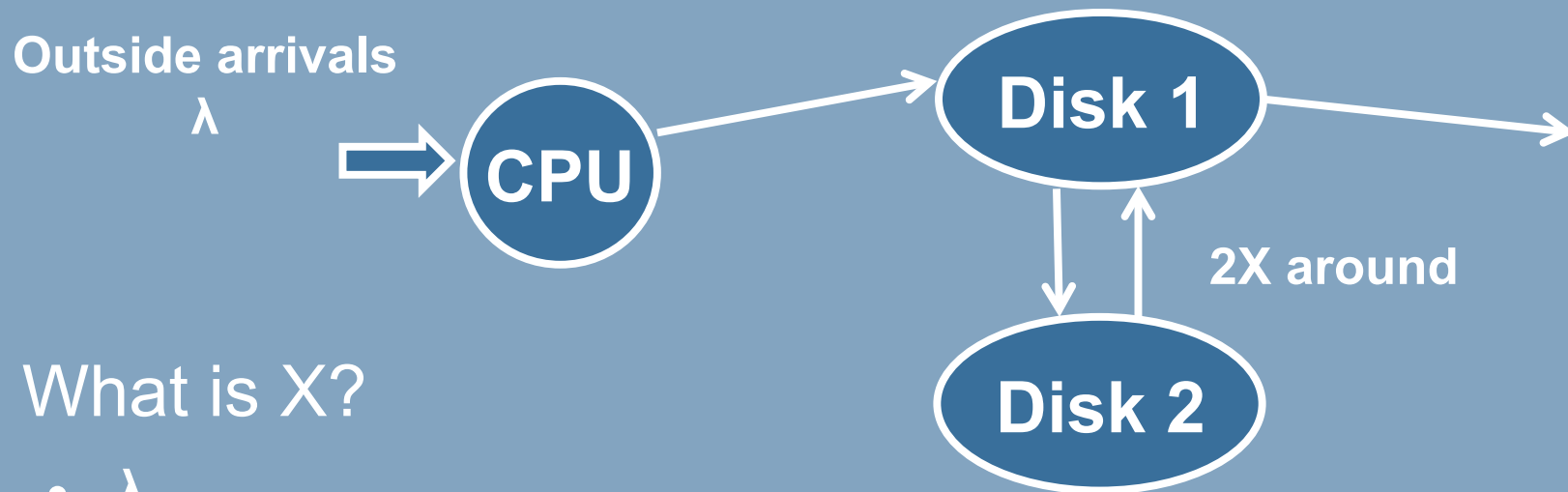
EXAMPLE 2: Network-of-queues with probabilistic routing

Outside arrivals



- What is X ?
 - $\sum r_i$
- What is X_i ?
 - $X_i = \lambda_i = r_i + \sum \lambda_j * P_{ji}$

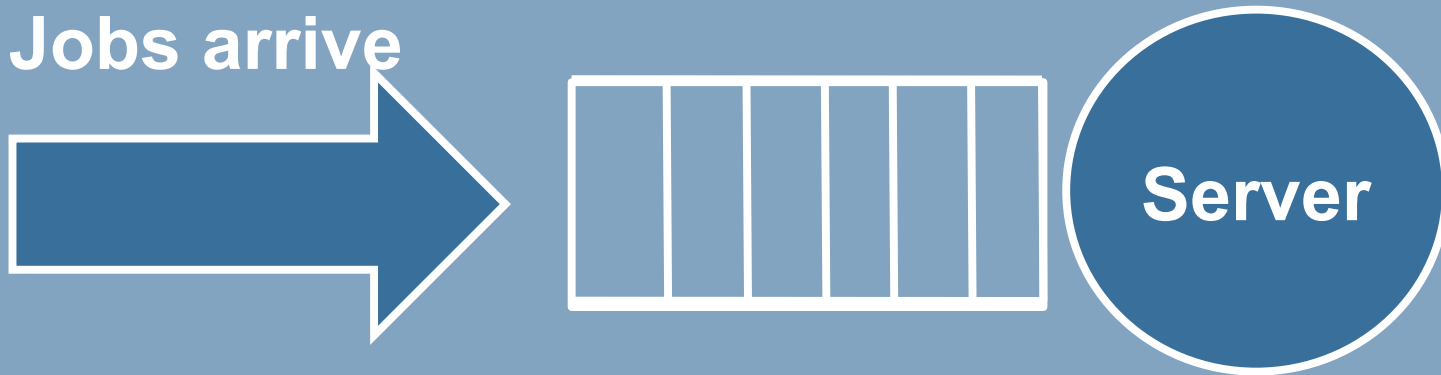
EXAMPLE 3: Network-of-queues with non-probabilistic routing



- What is X ?
 - λ
- What is X_{disk1}
 - 3λ
- What is X_{disk2}
 - 2λ

EXAMPLE 4: Finite buffer

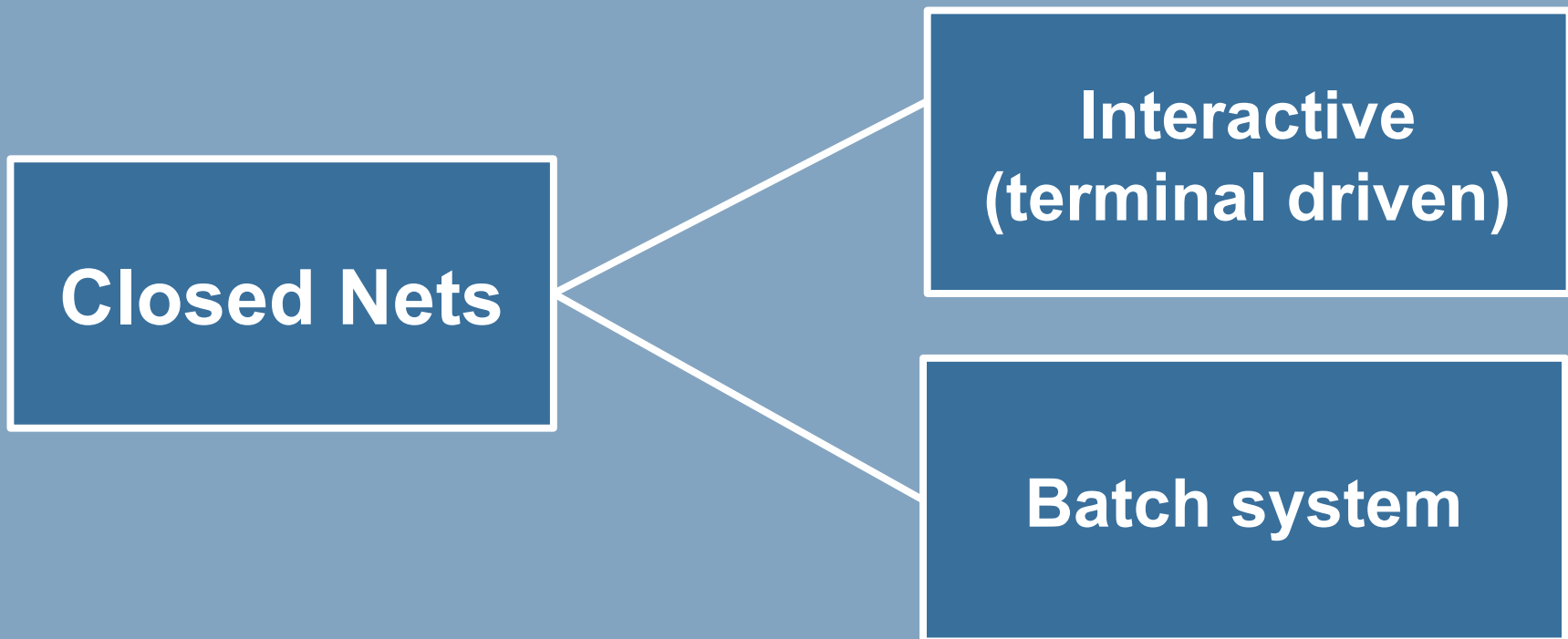
Jobs arrive



- Space in queue is limited to n jobs
- What is X ?
 - $X = U * \mu$
- But U is no longer λ / μ
- Need Markov chains ...

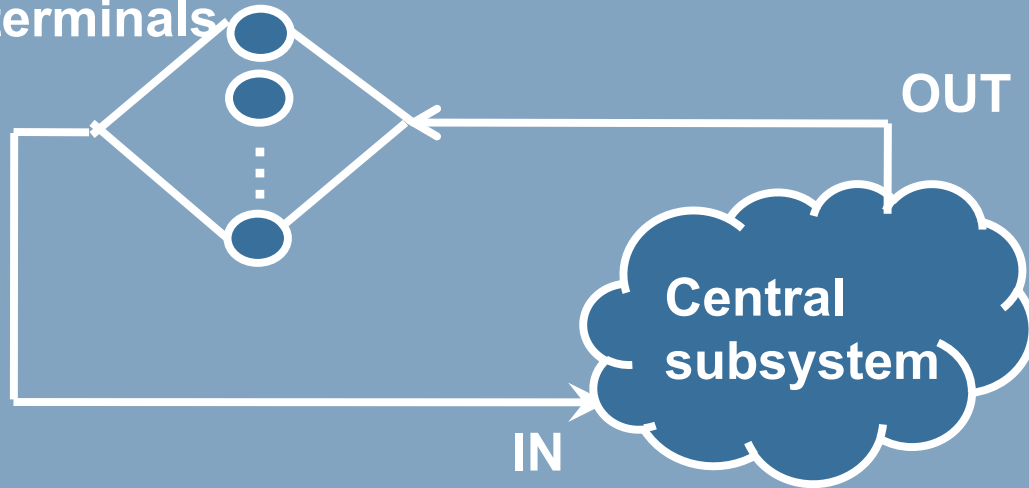
CLOSED SYSTEMS

- Closed networks have no external arrivals or departures



INTERACTIVE CLOSED SYSTEMS

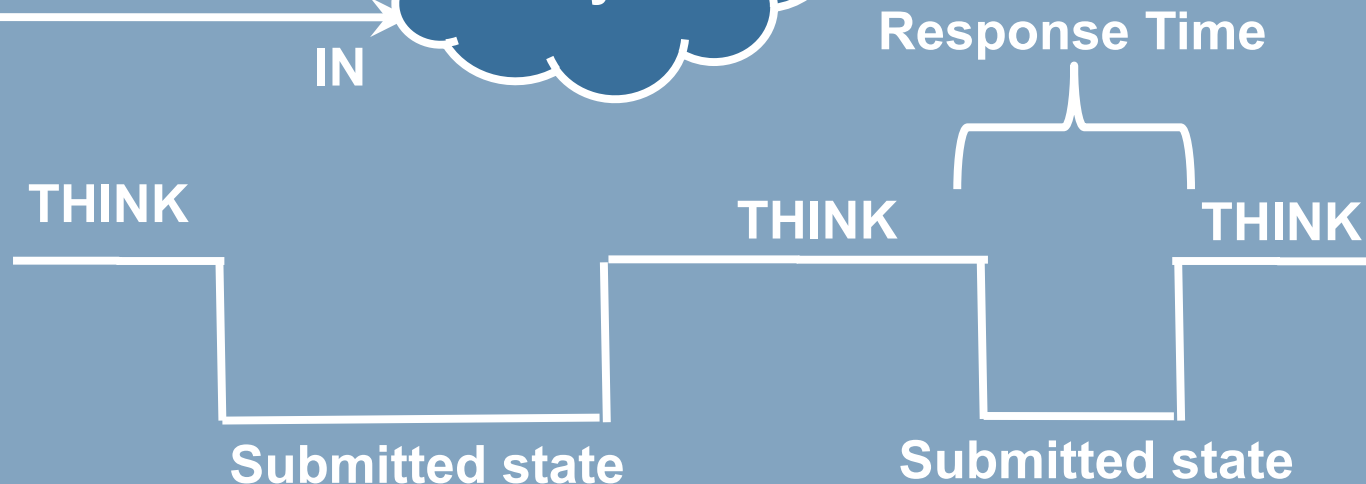
N user terminals



Parameters:

N = number of terminals

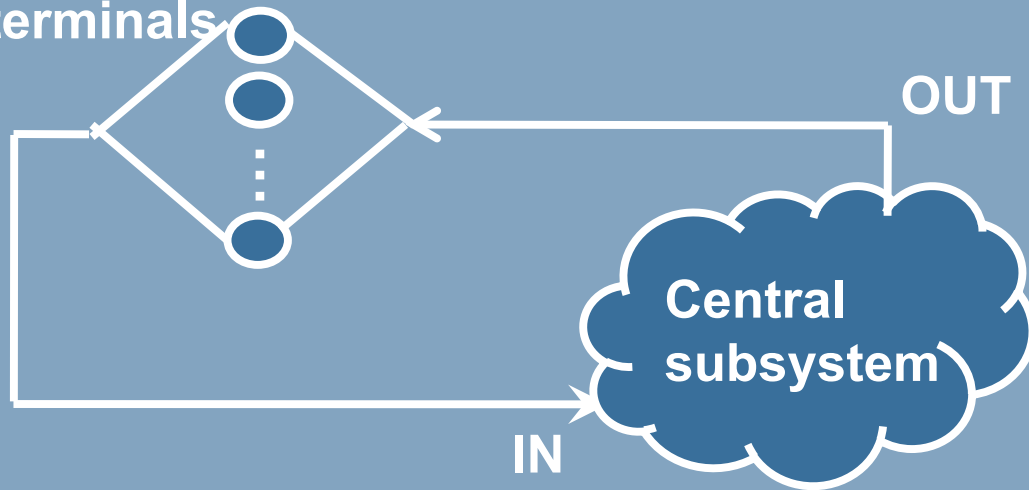
Think time Z



- How is response time defined?
 - Time it takes a job to go from IN to OUT

CLOSED INTERACTIVE SYSTEMS

N user
terminals



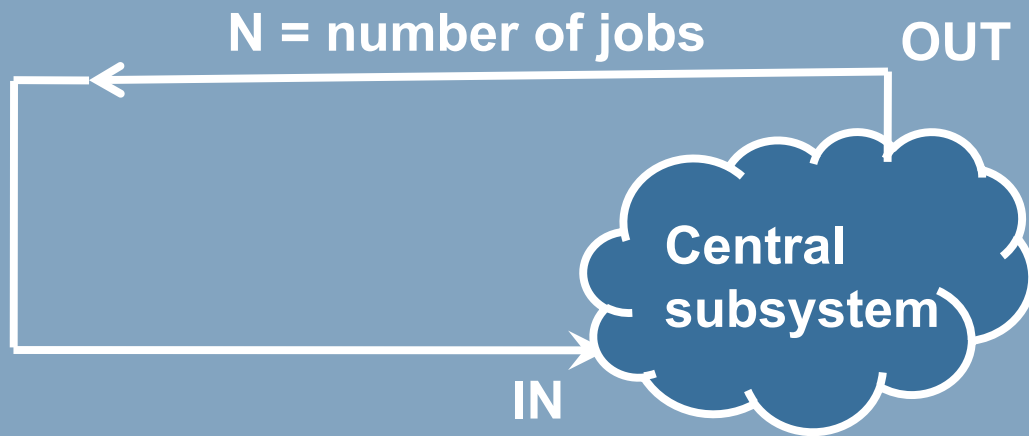
Parameters:

N = number of terminals

Think time Z

- Examples?
- Typical goals:
 - How high can we make N while keeping response times reasonably low?
 - Given a fixed N, what changes to central subsystem will improve response time the most?

CLOSED BATCH SYSTEMS

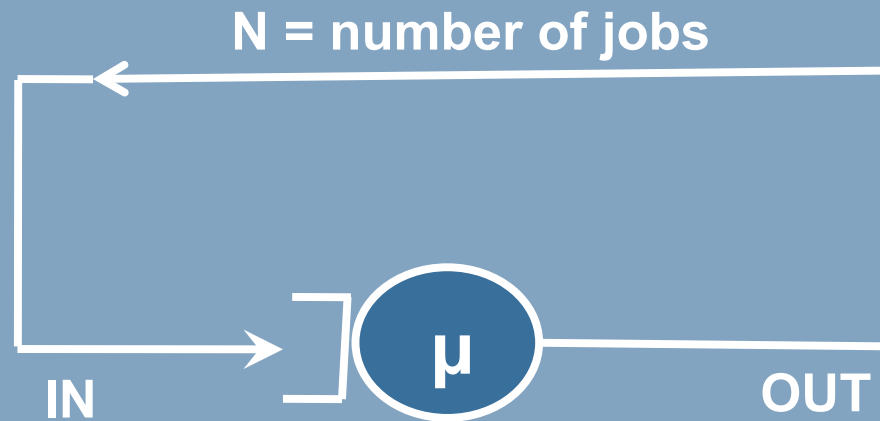


Parameters:

$N = \text{number of jobs}$

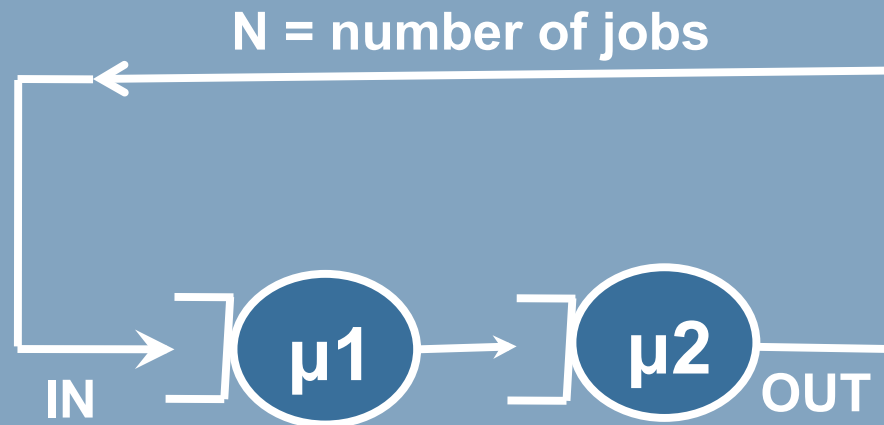
- As soon as one job completes, another one is started
- So there's always exactly N jobs in central subsystem
- Typical Goal?
 - Throughput!
- How is throughput defined?

THROUGHPUT IN BATCH SYSTEMS



- What is X ?
 - μ

THROUGHPUT IN BATCH SYSTEMS



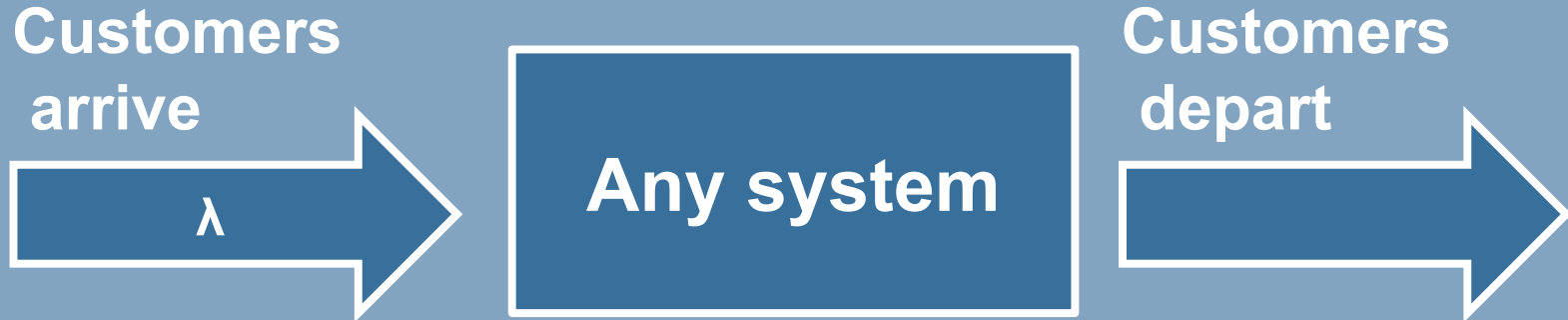
- What is X ?
 - It is NOT $\min \{\mu_1, \mu_2\}$
- Why?

SUMMARY SO FAR

- Open system
 - X independent of μ
 - X not affected by doubling μ
 - Throughput and response time independent
- Closed system
 - X depends on μ
 - Doubling μ affects X
 - Lower response time \Rightarrow higher X

LITTLE'S LAW

- The single most famous queueing result!
- For open systems:

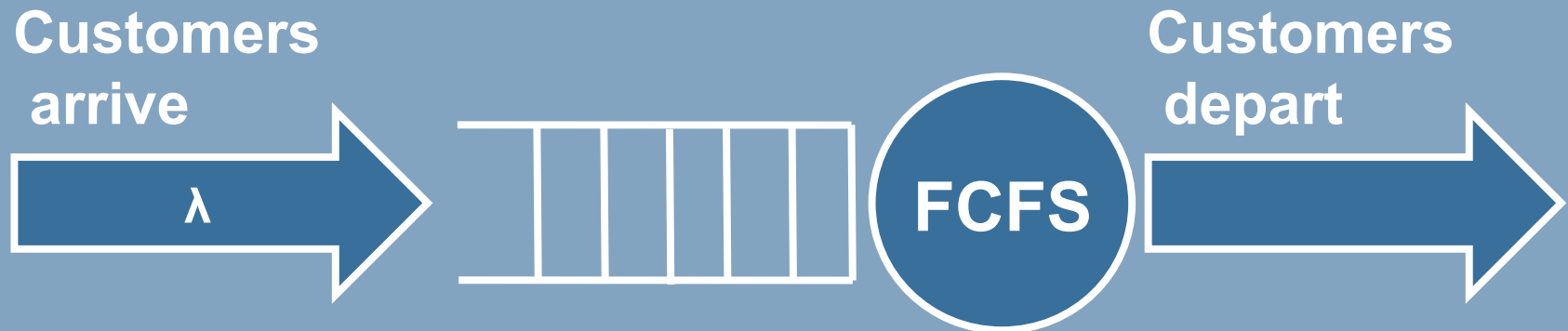


$$E[N] = \lambda * E[T]$$

Expected number
of jobs in system

Expected time in system/
Response time

LITTLE'S LAW: INTUITION



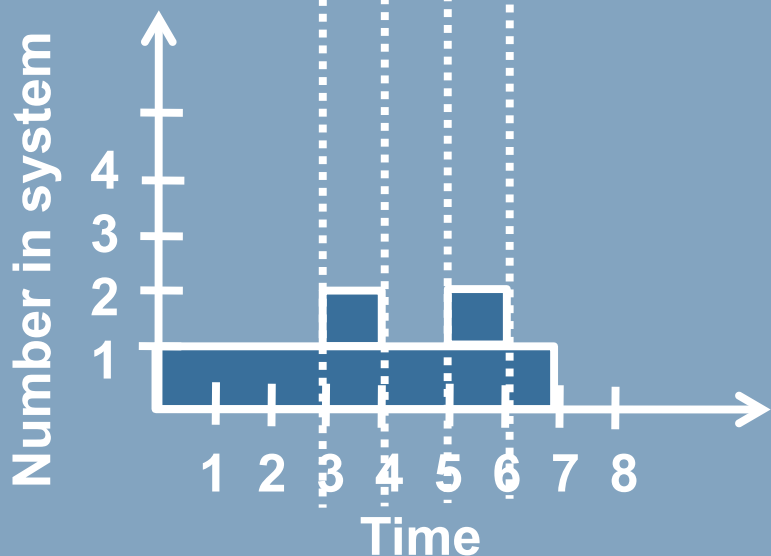
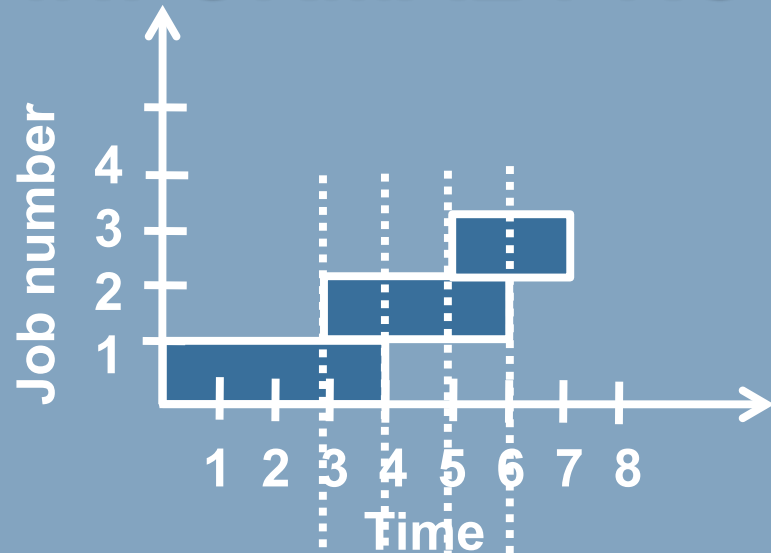
$$E[N] = \lambda * E[T]$$

- Customer arrives, sees $E[N]$ in system
- Average rate of completions is λ
- Expected time until customer leaves:

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

LITTLE'S LAW: INFORMAL PROOF

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

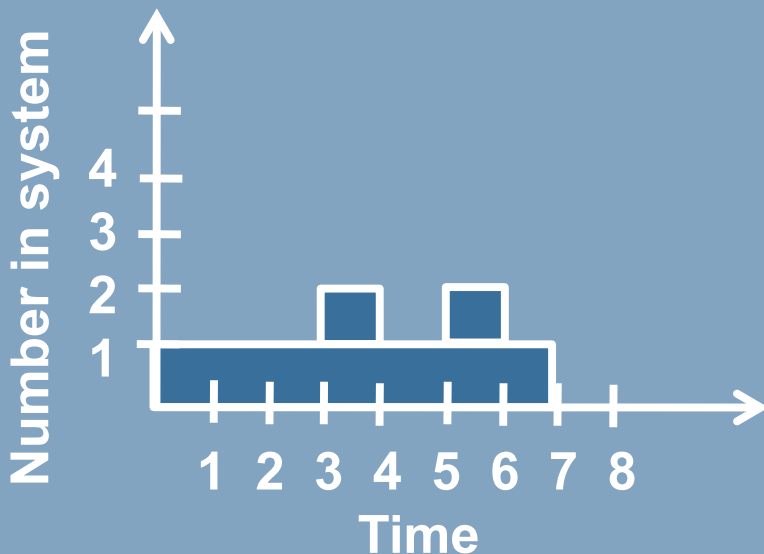
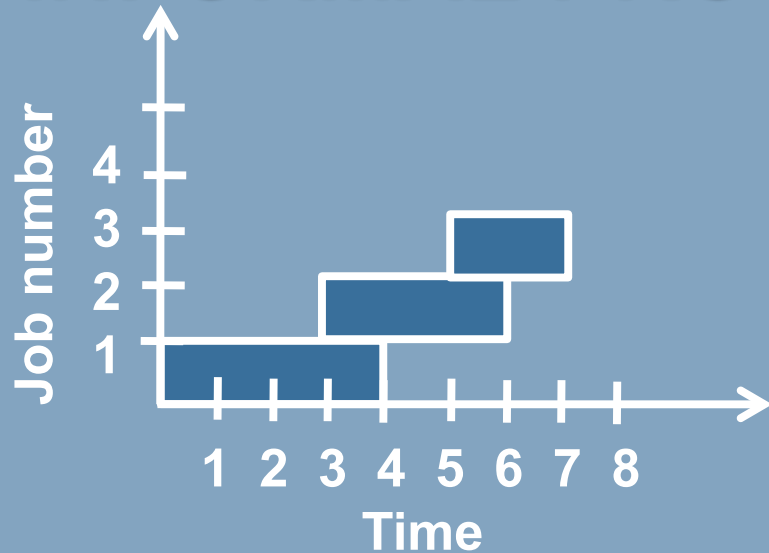


- How does the blue area in graph 1 compare to the blue area in graph 2?

=> they are identical!

LITTLE'S LAW: INFORMAL PROOF

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



- J = blue area
- n = number jobs completed by time t

- $E[T] = J/n$
- $E[N] = J/t$
 $= E[T] * n / t$
 $= E[T] * \lambda$

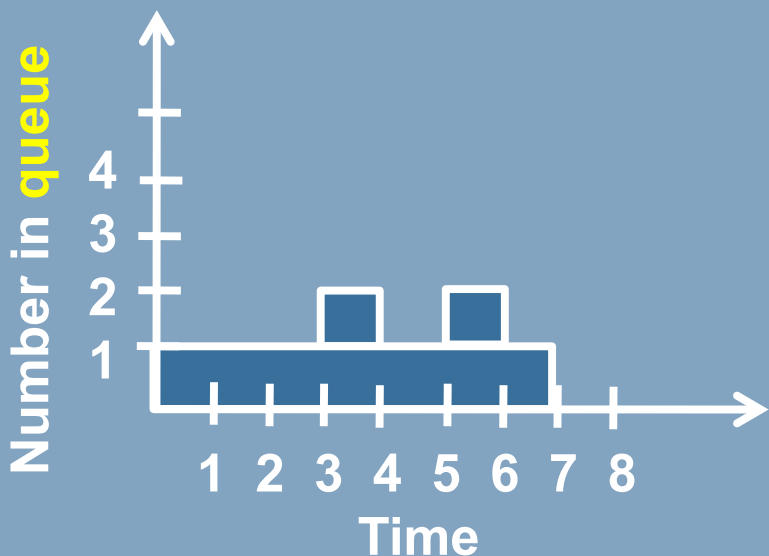
LITTLE'S LAW FOR TIME IN QUEUE



$$E[N_Q] = \lambda * E[T_Q]$$

LITTLE'S LAW FOR TIM IN QUEUE

$$E[N_Q] = \lambda * E[T_Q]$$



- J = blue area
- n = number jobs completed by time t

- $E[T_Q] = J/n$
- $E[N_Q] = J/t$
 $= E[T_Q] * n / t$
 $= E[T_Q] * \lambda$

LITTLE'S LAW FOR JOB CLASSES

- Suppose we're interested only in one type of jobs, e.g. "red jobs". Can we apply Little's law to only red jobs?

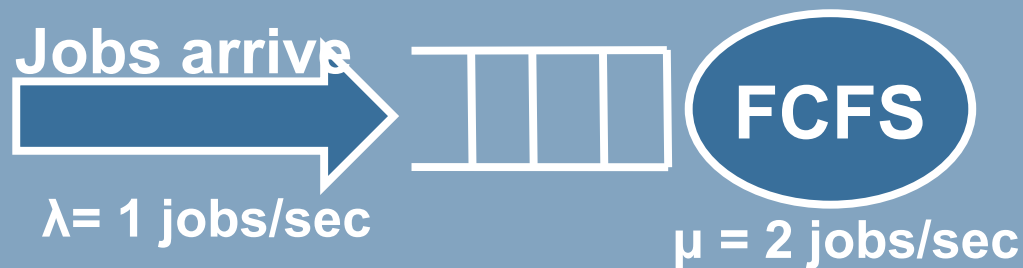
$$E[N_{\text{red}}] = \lambda_{\text{red}} * E[T_{\text{red}}]$$

APPLICATIONS OF LITTLE'S LAW?

Remember our pseudo proof of

$$U = \lambda/\mu \quad ?$$

PSEUDO PROOF OF $U = \lambda/\mu$

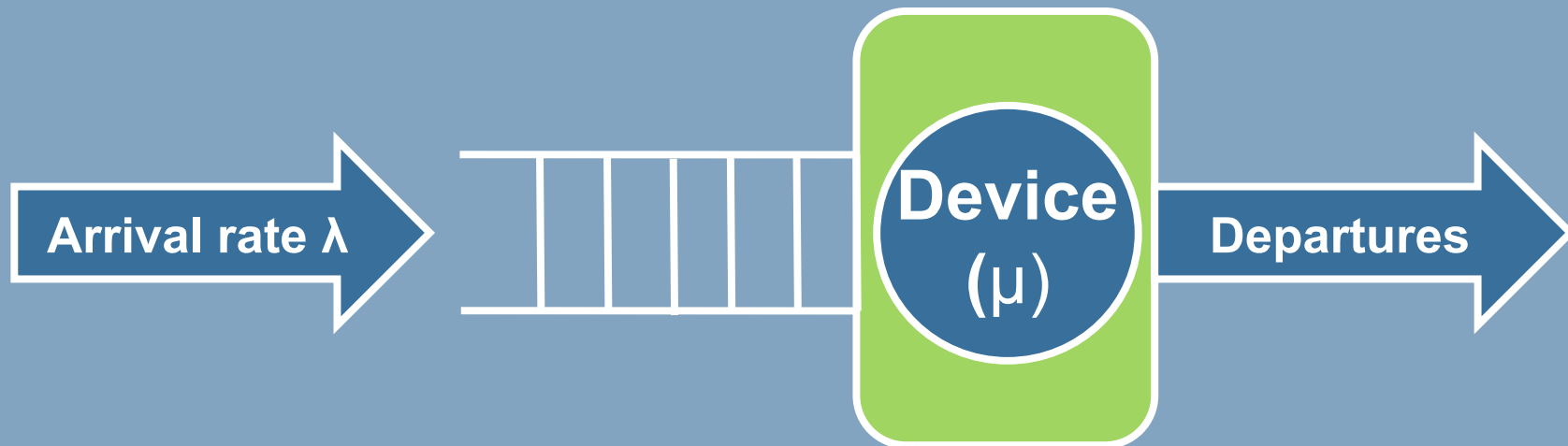


- What is U ?
 - Fraction of time server is busy
 - = mean service time / mean time betw. Arrivals
 - = $(1/\mu)/(1/\lambda) = \lambda/\mu$

APPLICATION OF LITTLE

$$\text{Little's Law: } E\{N\} = \lambda * E\{T\}$$

- Prove that $U = \lambda/\mu$!

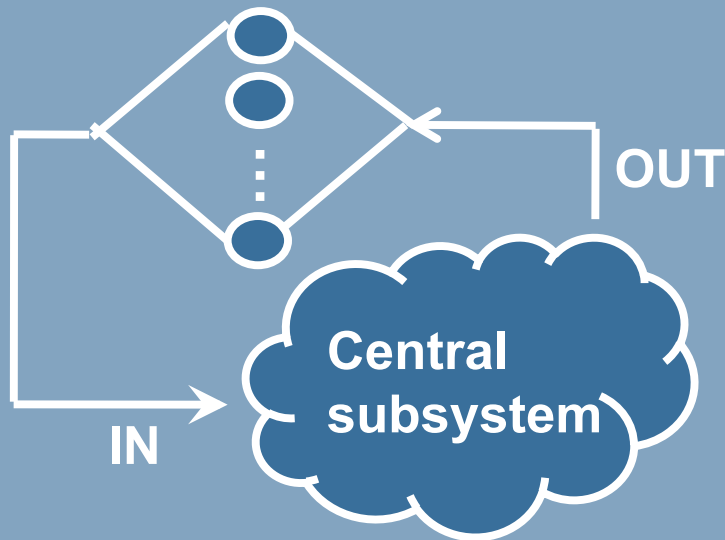


- Consider only the green subsystem
 - What is $E\{N\}$? $\Rightarrow U$
 - What is $E\{T\}$? $\Rightarrow E[S] = 1/\mu$
 - **Arrival rate?** $\Rightarrow \lambda$

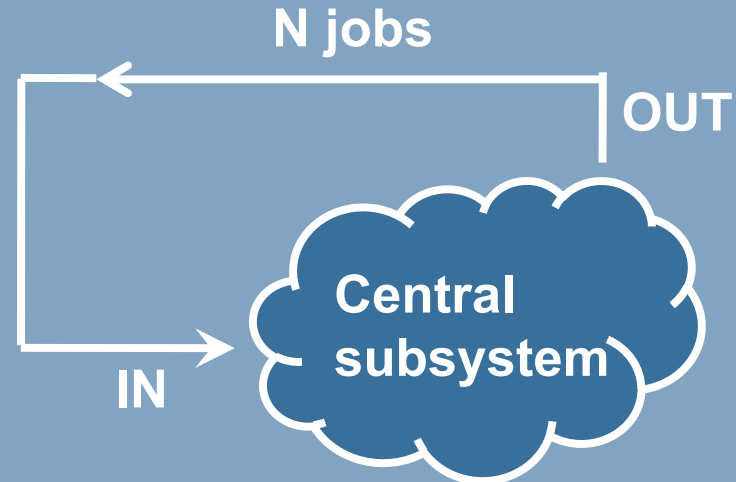
LITTLE'S LAW FOR CLOSED SYSTEMS

- Little's law in open systems: $E\{N\} = \lambda * E\{T\}$
- Little's law in closed systems? $N = X * E\{T\}$

N users (fixed) with thinktime Z



$E[T]$ = time from "OUT" to "OUT"
 $E[\text{RespTime}]$ = time from IN to "OUT"



$E[T]$ = time from "IN" to "OUT"

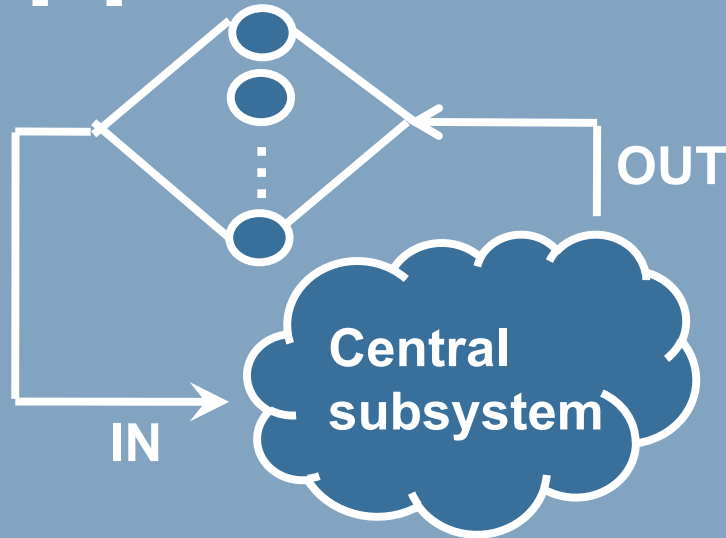
LITTLE'S LAW FOR CLOSED SYSTEMS

- Little's law in open systems: $E\{N\} = \lambda * E\{T\}$
- Little's law in closed systems? $N = X * E\{T\}$

- Main difference between open & closed?
 - In open, X does not depend on $E\{T\}$
 - In closed, X can be improved by improving $E\{T\}$

EXAMPLE 1

10 users
 $E[Z] = 5 \text{ sec}$



Little's Law

$$N = X * E\{T\}$$

$$\Rightarrow X = N/E[T]$$

$E[\text{RespTime}] = 15 \text{ sec}$

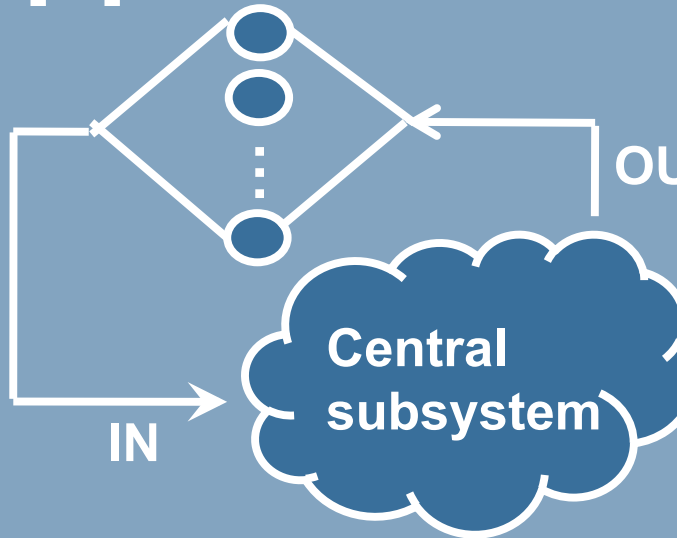
- What is the throughput of this system?
 - What is $E[T]$? $E[T] = E[Z] + E[\text{RespTime}]$
 $\Rightarrow X = N / (E[Z] + E[\text{RespTime}])$

Remember math from last week's paper:

In the flow-based *closed-loop model*, we have a fixed number of users N . Each user goes through a cycle of activity, with a flow of average size S followed by an idle period of average length T_i . The average flow arrival rate is given by $\lambda_c = N / (T_t + T_i)$, where T_t is the average flow completion time. The latter

EXAMPLE 1

10 users
 $E[Z] = 5 \text{ sec}$



Little's Law

Alternative version:
 $E[R] = N / X \text{ -- } E[Z]$
(Response Time Law)

$E[\text{RespTime}] = 15 \text{ sec}$

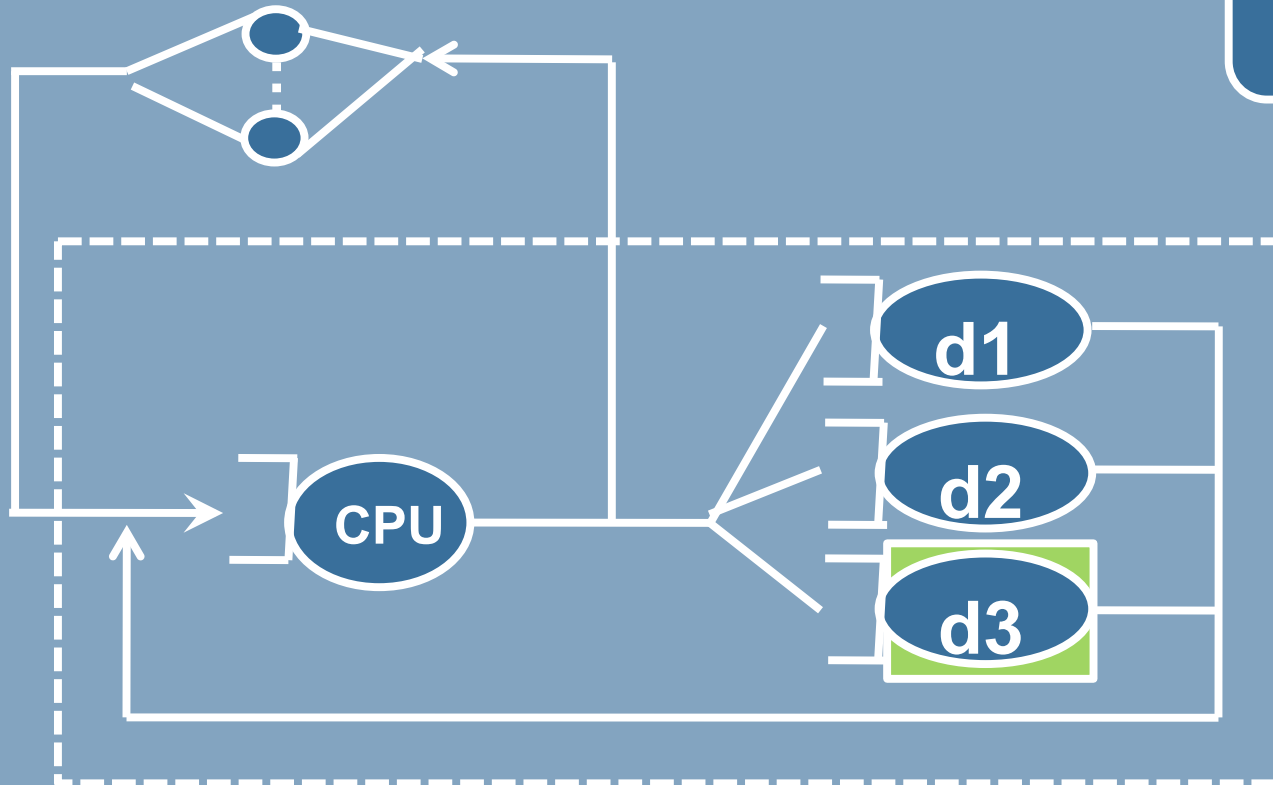
- What is the throughput of this system?
 - What is $E[T]$? $E[T] = E[Z] + E[\text{RespTime}]$
 $\Rightarrow X = N / (E[Z] + E[\text{RespTime}])$

EXAMPLE 2

Little's Law

$$N = X * E\{T\}$$

10 users



$$\begin{aligned} X_{d3} &= 40 \text{ jobs/sec} \\ E[S_{d3}] &= .0225 \text{ sec} \\ E[N_{d3}] &= 4 \end{aligned}$$

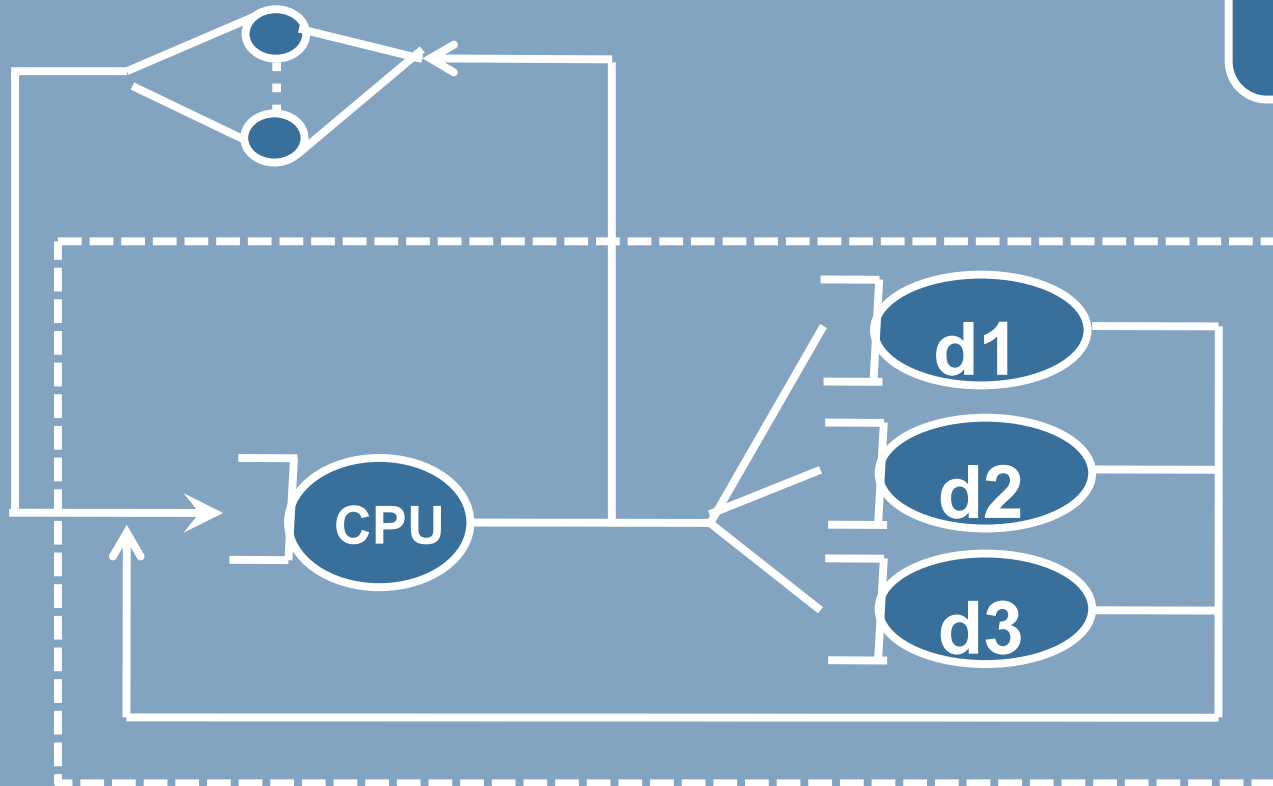
- What is the utilization of disk 3?
- $U_{d3} = E[\text{number of jobs in green box}]$
 $= X_{d3} * E[S_{d3}] = 90\%$

EXAMPLE 3

Little's Law

$$N = X * E\{T\}$$

10 users



$$\begin{aligned} X_{d3} &= 40 \text{ jobs/sec} \\ E[S_{d3}] &= .0225 \text{ sec} \\ E[N_{d3}] &= 4 \end{aligned}$$

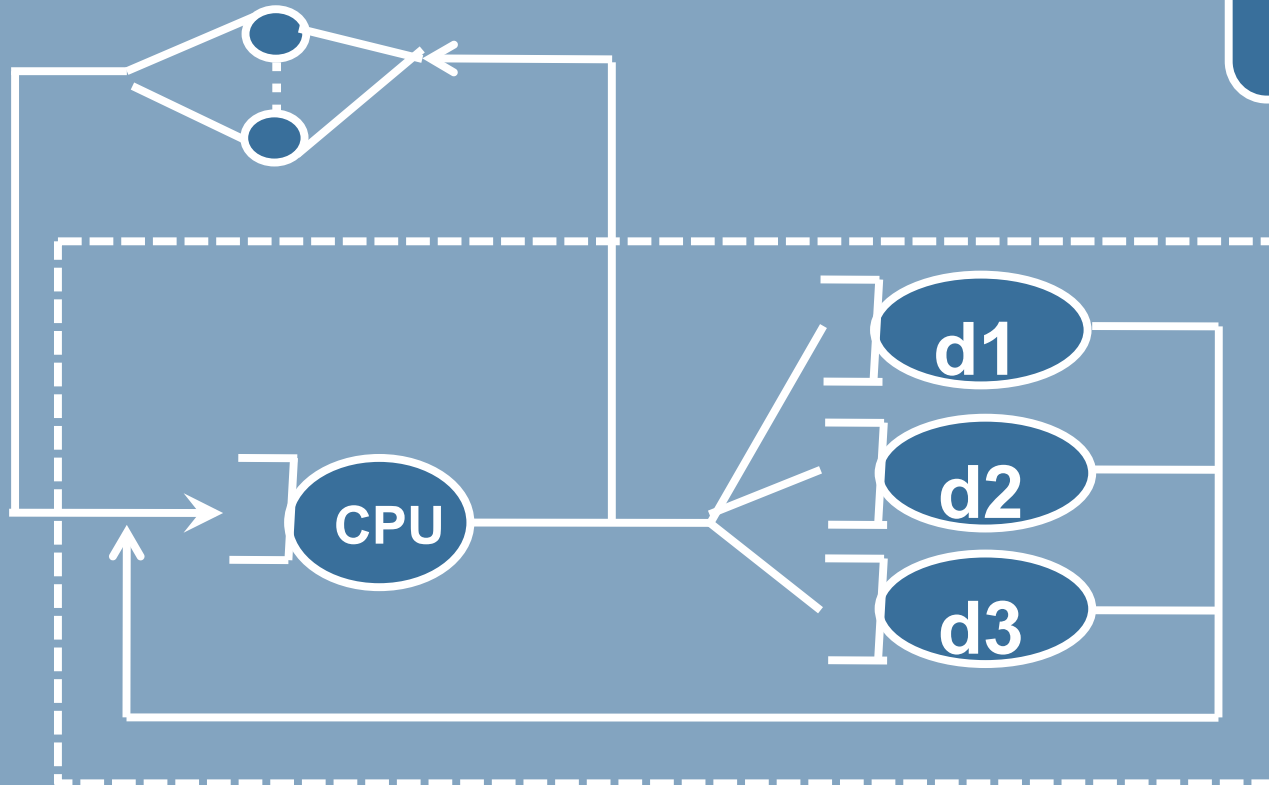
- What is T_Q of disk 3?
- $E[T_{d3}] = N_{d3} / X_{d3} = 4/40 = 0.1 \text{ sec}$
- $E[TQ_{d3}] = E[T_{d3}] - E[S_{d3}] = 0.0775 \text{ sec}$

EXAMPLE 4

Little's Law

$$N = X * E\{T\}$$

10 users



$$\begin{aligned} X_{d3} &= 40 \text{ jobs/sec} \\ E[S_{d3}] &= .0225 \text{ sec} \\ E[N_{d3}] &= 4 \end{aligned}$$

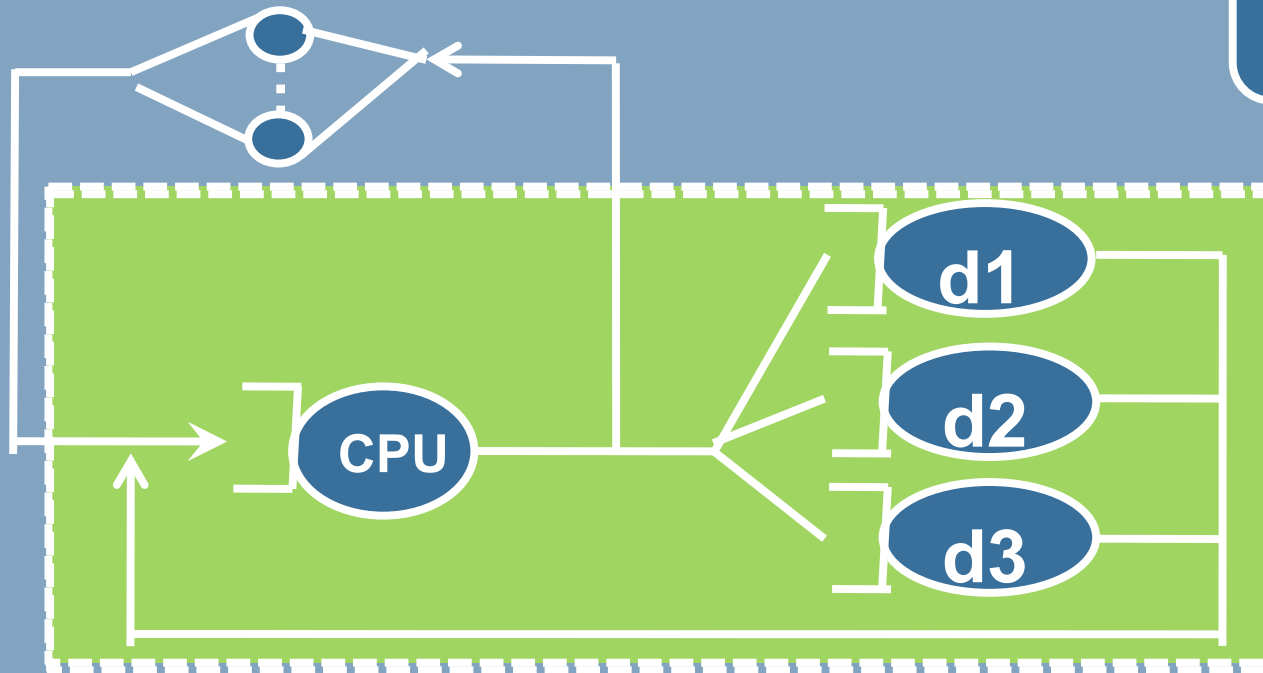
- What is $E[\text{number of requests queued at } d3]$?
- $E[NQ_{d3}] = E[N_{d3}] - E[\text{number in service}]$
 $= 4 - 0.9 = 3.1$

EXAMPLE 4

Little's Law

$$N = X * E\{T\}$$

10 users



$$E[\text{ready users}] = 7.5$$
$$E[Z] = 5 \text{ sec}$$

- What is the system throughput ?
 - $X = N / E[T] = N / (E[Z] + E[R]) = 10 / (5 + E[R])$
- But what is $E[R]$?
 - $E[R] = E[N_{\text{greenbox}}] / X_{\text{greenbox}} = 7.5 / X$

REVIEW

- How to relate the expected number of jobs in the system to expected time spent in the system?

- Open system:

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

- Closed batch system
(zero thinktime)

$$N = X * E[T]$$

- Closed interactive system

$$E[R] = N/X - E[Z]$$

REVIEW: UTILIZATION LAW

- How can you relate device utilization and device throughput?

$$U_i = X_i * E[S_i]$$

- How can you prove this? Apply Little's law to individual device

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