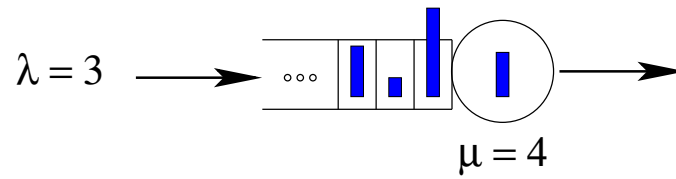## A single-server queue



$\lambda = 3$

$\mu = 4$

**Average Arrival Rate:**

**Interarrival Time:**

**Mean Interarrival Time:**

**Job Size (Service Requirement):**

**Mean Job Size:**

**Average Service Rate:**

$\lambda = 3$ ⟶ 

$\mu = 4$
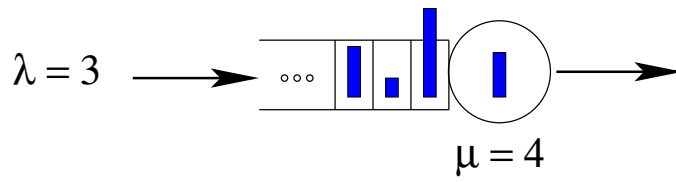
**PEOPLE SPEAK** vs. **QUEUEING SPEAK**

**Examples from your work:**

# Common Performance Metrics

- Response Time, $T$:

- Waiting Time or Delay, $T_Q$

- Number of jobs in system, $N$

- Number of jobs in queue, $N_Q$

# Stability



$\lambda = 3$    $\mu = 4$

**Question:** What happens if $\lambda > \mu$?

Here's why:
$N(t) =$ Number jobs in system *at* time $t$
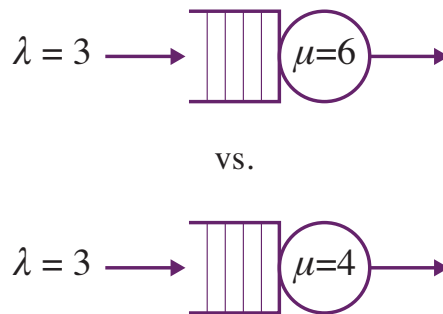$A(t) =$ Number arrivals *by* time $t$
$C(t) =$ Number completions(depatures) *by* time $t$
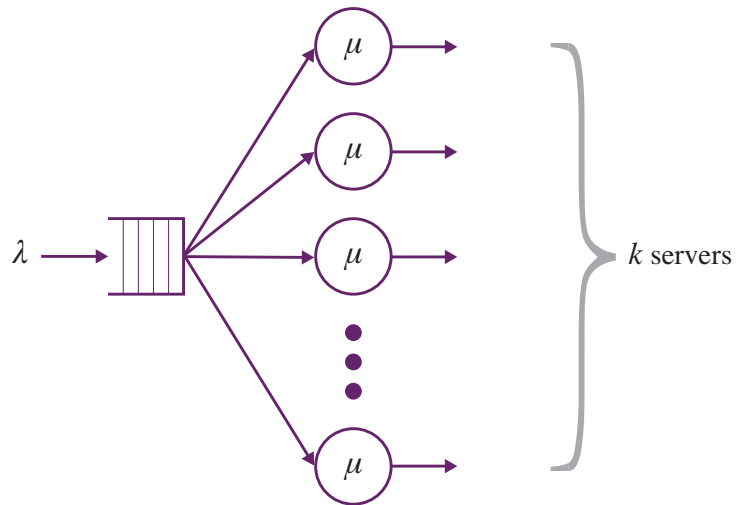
**We will always assume $\lambda < \mu$. (Stability)**

# Throughput, X

**Question:** What is throughput?

**Question:** Which has higher throughput?

$$\lambda = 3 \longrightarrow \boxed{\phantom{||}} \, (\mu = 6) \longrightarrow$$

vs.

$$\lambda = 3 \longrightarrow \boxed{\phantom{||}} \, (\mu = 4) \longrightarrow$$
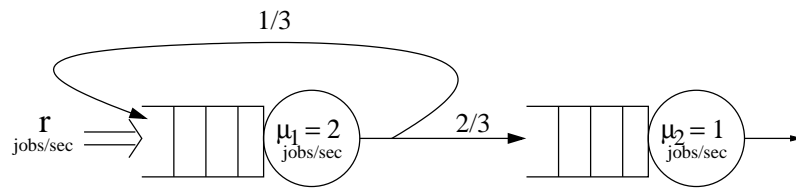
# Throughput for Server Farm



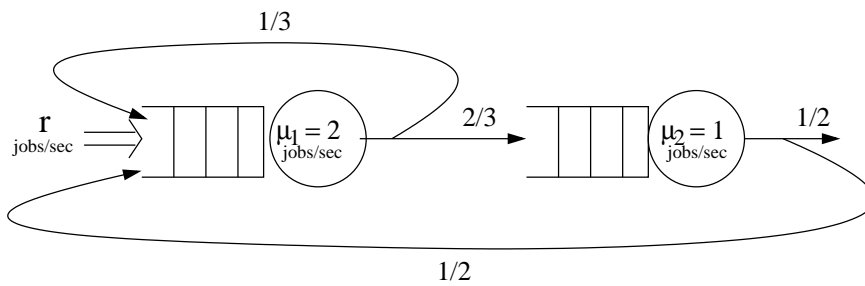**Question:** What condition is needed for stability?

**Question:** What's the throughput? (assuming stability)

# Throughput for Network of Queues

**Question:** What is the maximum outside arrival rate?
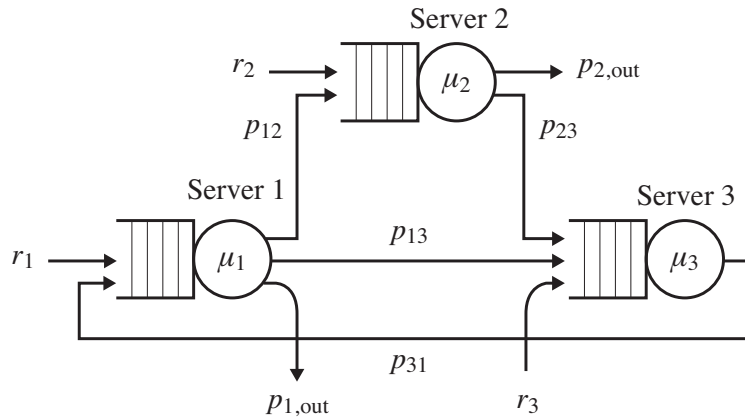


(a)
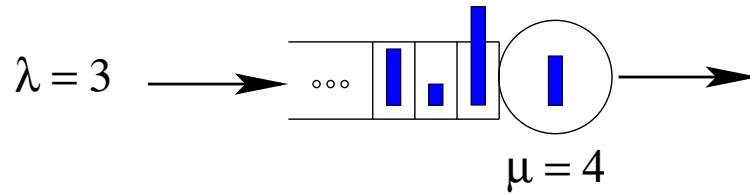


(b)

# Throughput for Network of Queues



**Question:** What's the throughput of this system? (assume stability)

**Question:** What's the throughput of server $i$?

**Question:** What do we need for stability of this system?

# Device Utilization (Load)

When talking about "utilization," we're thinking of a single device.



$\lambda = 3$      $\mu = 4$

Defn: Device **utilization**, a.k.a. **load**, is the long-run fraction of time that the device is busy.

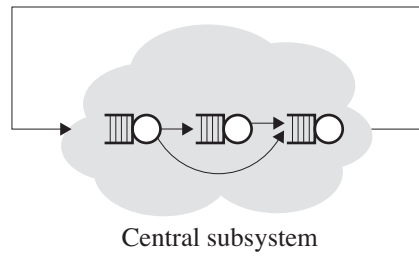Let $B(t) =$ total time server is busy during $[0, t]$.
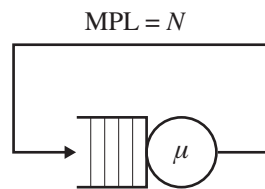Q: What is $\rho$?

Example:
3 jobs/sec arrive on average.
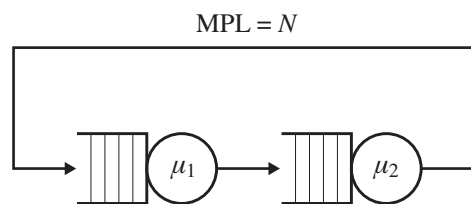Each job contributes 1/4 sec of work on average.
Q: What is $\rho$?

# Closed System (Batch)



Central subsystem

**Question:** What is throughput below?



MPL = $N$

**Question:** What is throughput below?



MPL = $N$

# Closed System (Interactive)



*N* user terminals

in

out

Central subsystem



Think state

Think state

Think state

Submitted state

Submitted state

in

out