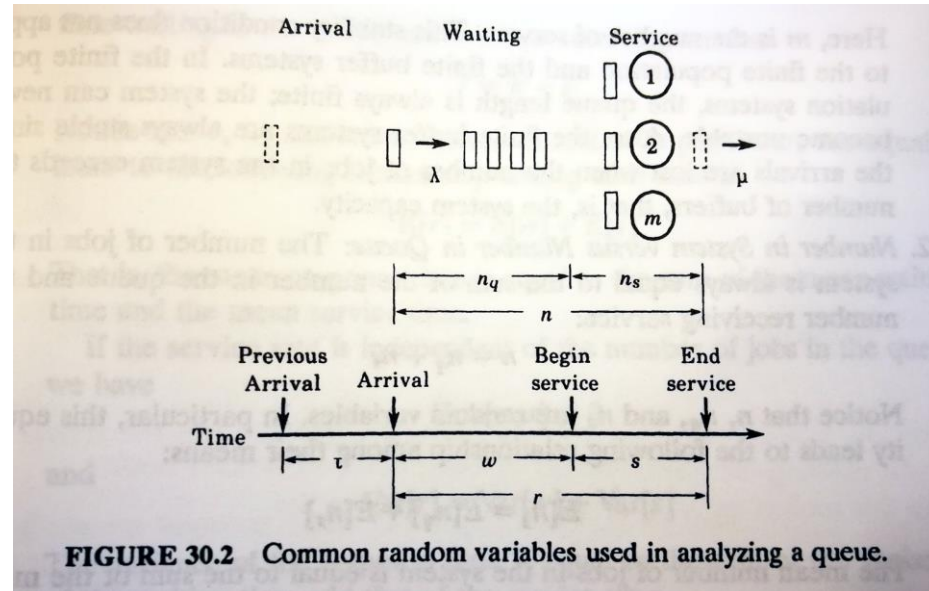


# Queueing Theory

M/M/1 and M/M/m Queues

# Why queuing theory?

- Reason about the system
- Analyze different components
- Determine bottlenecks
- Queues are a core aspect of complex systems



# Little's Law

**Number of requests in the system = arrival rate \* mean response time**

- Relates the number of requests to the response time
- Valid for any type of queueing system
- Valid for systems in its entirety or for parts of the system

**Number of requests in the queue = arrival rate \* mean waiting time in the queue**

# What you are trying to predict?

- Be clear what the input to the model is. This depends on what you have measured!
- Be clear what elements you want to predict
- Explain which aspects of your system correspond to the model and which ones do not
- Consider all the data you have gathered  
There are many options to proceed!



Input	Anything that has been measured
Goal	Arrival rate, service rate
Output	Length of the queues, number of request in the system, graphs, behaviour, ...

**Box 31.1 M/M/1 Queue**

- Parameters:  
 $\lambda$  = arrival rate in jobs per unit time  
 $\mu$  = service rate in jobs per unit time
- Traffic intensity:  $\rho = \lambda/\mu$
- Stability condition: Traffic intensity  $\rho$  must be less than 1.
- Probability of zero jobs in the system:  $p_0 = 1 - \rho$
- Probability of  $n$  jobs in the system:  $p_n = (1 - \rho)\rho^n$ ,  $n = 0, 1, \dots, \infty$
- Mean number of jobs in the system:  $E[n] = \rho/(1 - \rho)$
- Variance of number of jobs in the system:  $\text{Var}[n] = \rho/(1 - \rho)^2$
- Probability of  $k$  jobs in the queue:

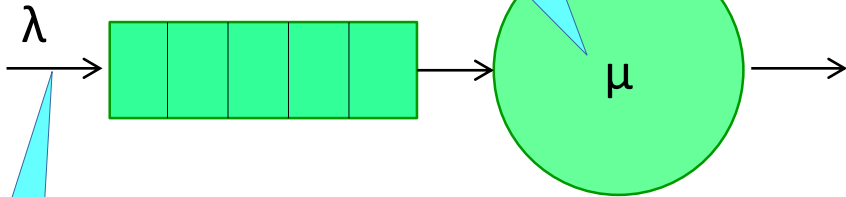
$$P(n_q = k) = \begin{cases} 1 - \rho^2, & k = 0 \\ (1 - \rho)\rho^{k+1}, & k > 0 \end{cases}$$

- Mean number of jobs in the queue:  $E[n_q] = \rho^2/(1 - \rho)$
- Variance of number of jobs in the queue:  
 $\text{Var}[n_q] = \rho^2(1 + \rho - \rho^2)/(1 - \rho)^2$
- Cumulative distribution function of the response time:  
 $F(r) = 1 - e^{-r\mu(1-\rho)}$
- Mean response time:  $E[r] = (1/\mu)/(1 - \rho)$
- Variance of the response time:  $\text{Var}[r] = \frac{1/\mu^2}{(1 - \rho)^2}$
- $q$ -Percentile of the response time:  $E[r] \ln[100/(100 - q)]$
- 90-Percentile of the response time:  $2.3E[r]$
- Cumulative distribution function of waiting time:  
 $F(w) = 1 - \rho e^{-\mu w(1-\rho)}$
- Mean waiting time:  $E[w] = \rho \frac{1/\mu}{1 - \rho}$
- Variance of the waiting time:  $\text{Var}[w] = (2 - \rho)\rho/[\mu^2(1 - \rho)^2]$
- $q$ -Percentile of the waiting time:  $\max\left(0, \frac{E[w]}{\rho} \ln[100\rho/(100 - q)]\right)$
- 90-Percentile of the waiting time:  $\max\left(0, \frac{E[w]}{\rho} \ln[10\rho]\right)$
- Probability of finding  $n$  or more jobs in the system:  $\rho^n$
- Probability of serving  $n$  jobs in one busy period:  

$$\frac{1}{n} \binom{2n-2}{n-1} \frac{\rho^{n-1}}{(1 + \rho)^{2n-1}}$$

# M/M/1 Queue

How can we determine this?



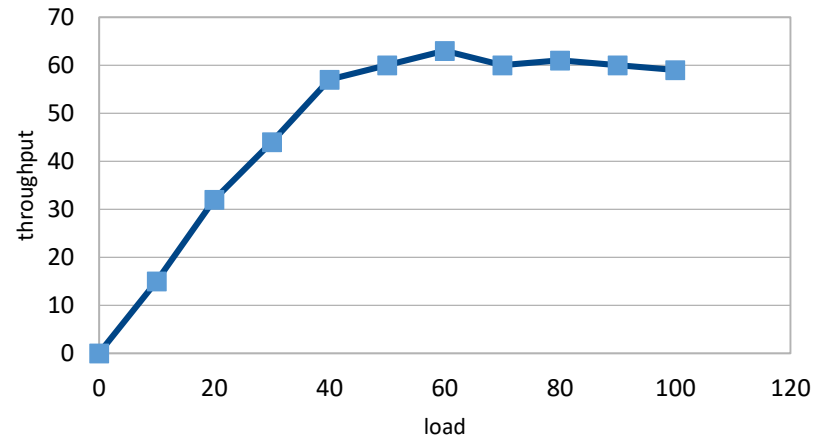
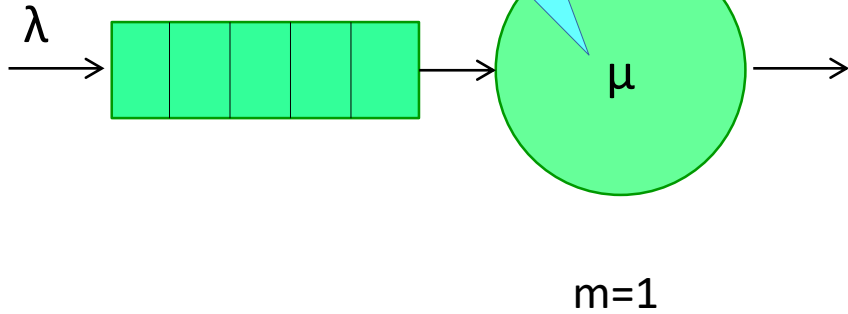
Measured

$m=1$

Given by the model

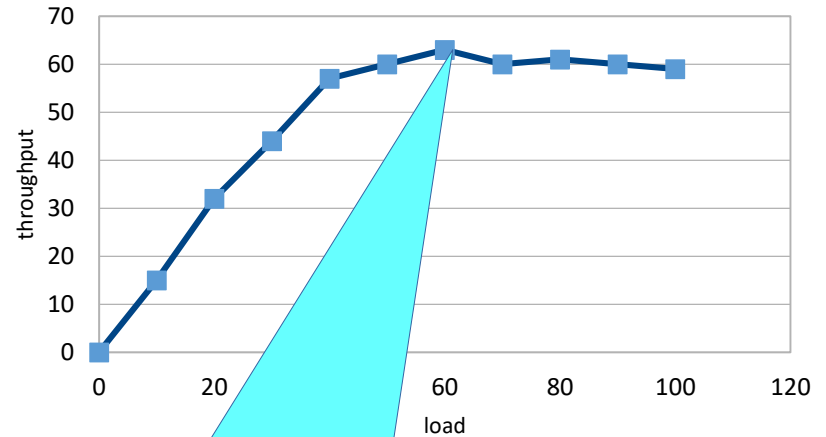
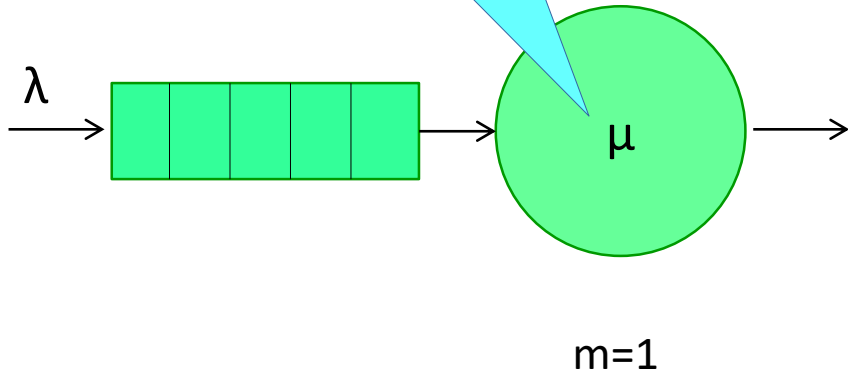
# M/M/1 Queue

How can we determine this?



# M/M/1 Queue

How can we determine this?



Find the maximum observed throughput

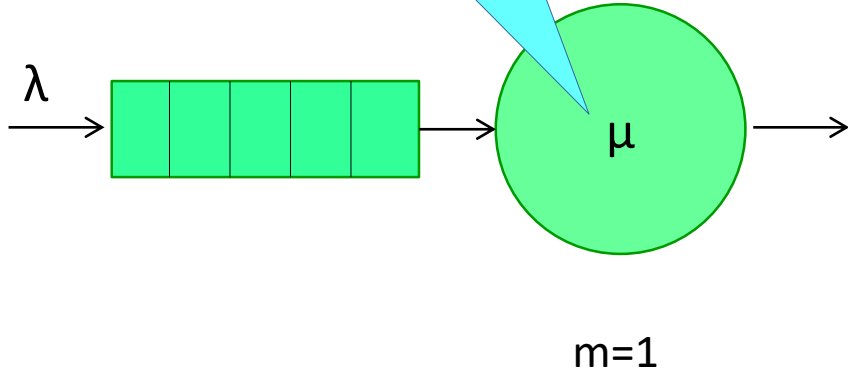
- Input to the model:  $\lambda$ ,  $\mu$
- How to measure the service rate  $\mu$ ?

There are many approaches, depending what aspect of your system you want to model!

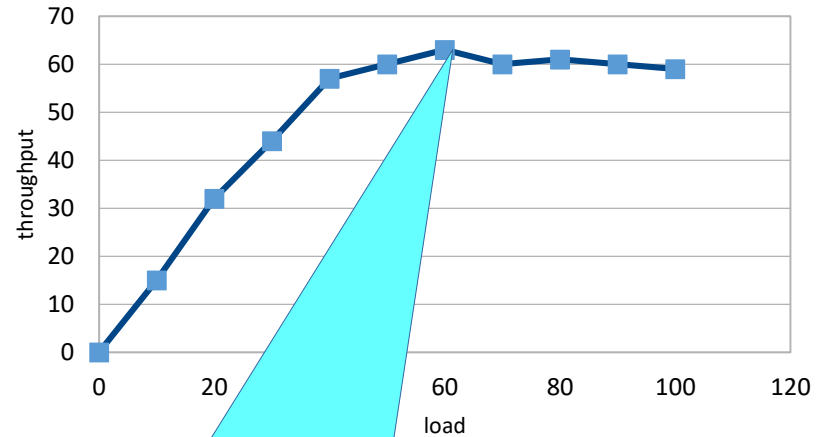
E.g., use configuration that gives maximum throughput

# M/M/1 Queue

How can we determine this?



Service rate > arrival rate.  $\rho < 1$ .



Find the maximum observed throughput

• Input to the model:  $\lambda$ ,  $\mu$

• How to measure the service rate  $\mu$ ?

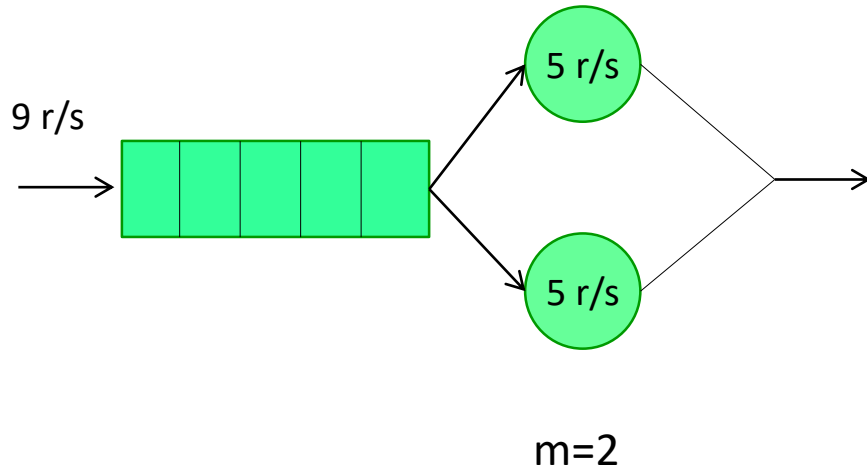
– There are many approaches, depending what aspect of your system you want to model!

– E.g., use configuration that gives maximum throughput



Example of computing the output of  
M/M/m models

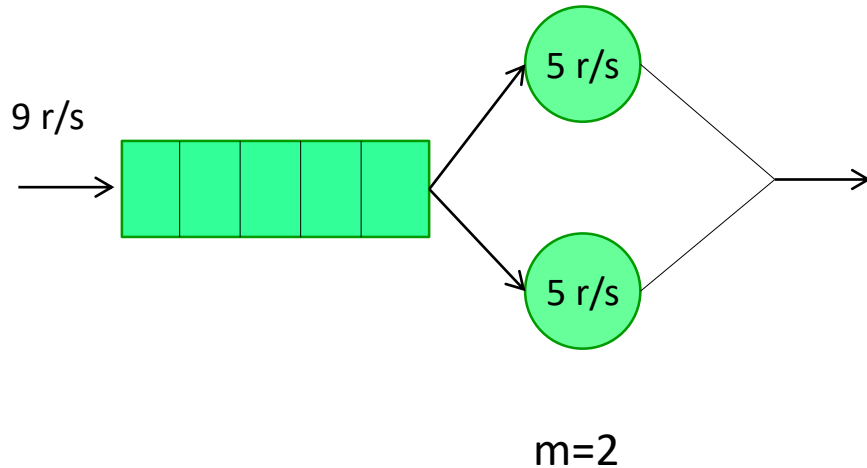
# Example 1: M/M/2



**Input:  $\lambda=9$  req/s,  $m=2$ ,  $\mu=5$  req/s**

- Utilization:
- Average service time per worker:
- Average number of requests in queue:
- Average waiting time in queue:
- Average response time:

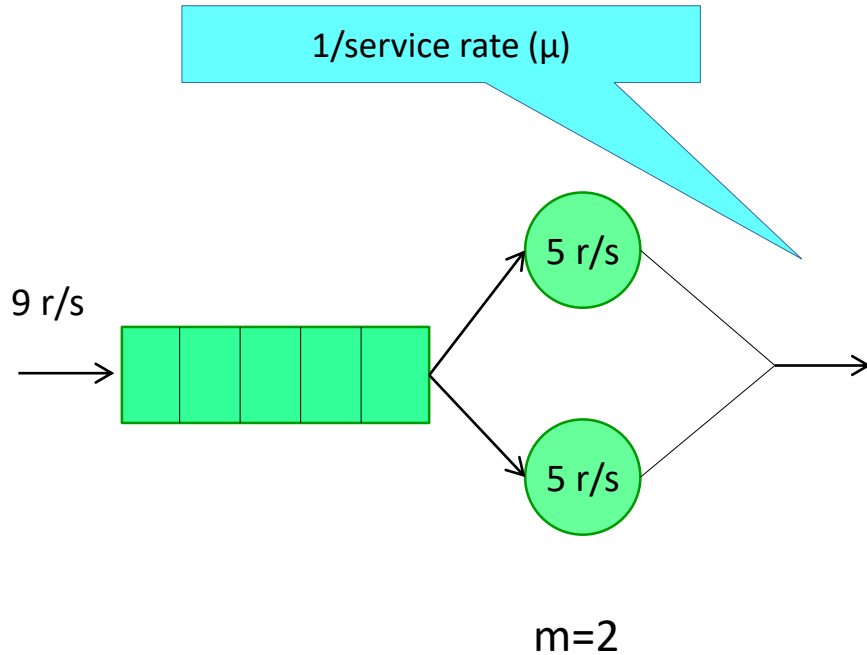
# Example 1: M/M/2



**Input:  $\lambda=9 \text{ req/s}$ ,  $m=2$ ,  $\mu=5 \text{ req/s}$**

- Utilization:  $\rho = \lambda / (m * \mu) = 90\%$
- Average service time per worker:
- Average number of requests in queue:
- Average waiting time in queue:
- Average response time:

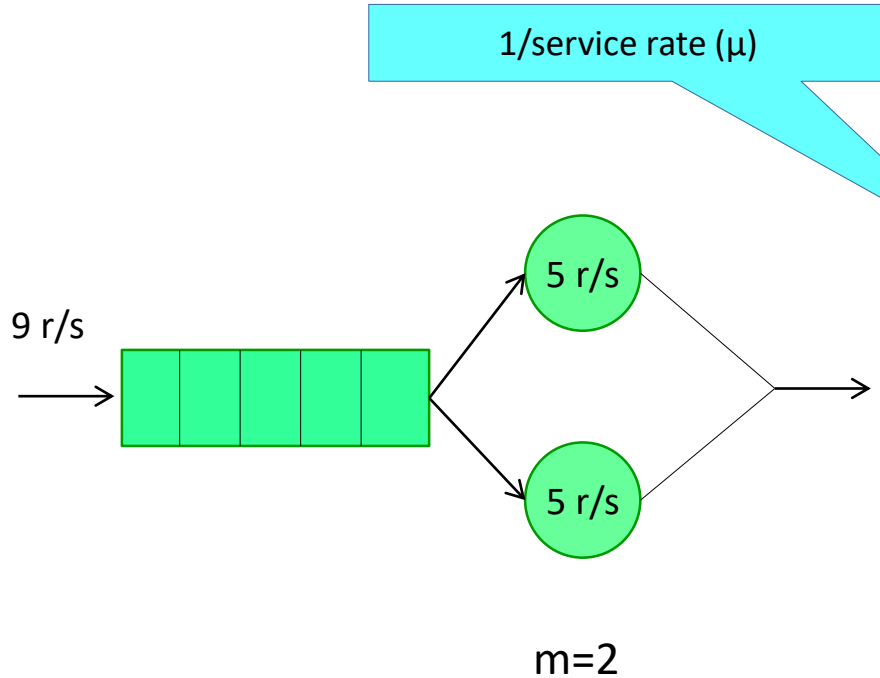
# Example 1: M/M/2



**Input:  $\lambda=9 \text{ req/s}$ ,  $m=2$ ,  $\mu=5 \text{ req/s}$**

- Utilization:  $\rho = \lambda / (m * \mu) = 90\%$
- Average service time per worker:  
 $E[s] = 0.20 \text{ s}$
- Average number of requests in queue:  
 $E[nq] = 7.674$
- Average waiting time in queue:
- Average response time:

# Example 1: M/M/2



Input:  $\lambda=9$  req/s,  $m=2$ ,  $\mu=5$  req/s

- Utilization:  $\rho = \lambda / (m * \mu) = 90\%$
- Average service time per worker:  
 $E[s] = 0.20$  s
- Average number of requests in queue:  
 $E[nq] = 7.674$
- Average waiting time in queue:  
 $E[wq] = 0.853$  s
- Average response time:

Mean number of req. in the system  
= arrival rate \* mean response time  
For the queue:  $E[nq] = \lambda * E[wq]$   
**Little's Law!**

# Example 1: M/M/2

1/service rate ( $\mu$ )

Mean number of req. in the system  
= arrival rate \* mean response time  
For the queue:  $E[nq] = \lambda * E[wq]$   
**Little's Law!**

Total response time  
= queuing time + service time

Can also be computed as a function  
of  $\lambda, \mu, m, \rho$ .. (see book)

**Input:  $\lambda=9$  req/s,  $m=2$ ,  $\mu=5$  req/s**

- Utilization:  $\rho = \lambda / (m * \mu) = 90\%$
- Average service time per worker:  
 $E[s] = 0.20$  s
- Average number of requests in queue:  
 $E[nq] = 7.674$
- Average waiting time in queue:  
 $E[wq] = 0.853$  s
- Average response time:  
 $E[w] = 0.853 + 0.20 = 1.053$

# Example 1: M/M/2

1/service rate ( $\mu$ )

Mean number of req. in the system  
= arrival rate \* mean response time  
For the queue:  $E[nq] = \lambda * E[wq]$   
**Little's Law!**

Total response time  
= queuing time + service time

Can also be computed as a function  
of  $\lambda$ ,  $\mu$ ,  $m$ ,  $\rho$ .. (see book)

**Input:  $\lambda=9$  req/s,  $m=2$ ,  $\mu=5$  req/s**

• Utilization:  $\rho = \lambda / (m * \mu) = 90\%$

• Average service time per worker:

$$E[s] = 0.20 \text{ s}$$

• Average number of requests in queue:

$$E[nq] = 7.67$$

• Average waiting time in queue:

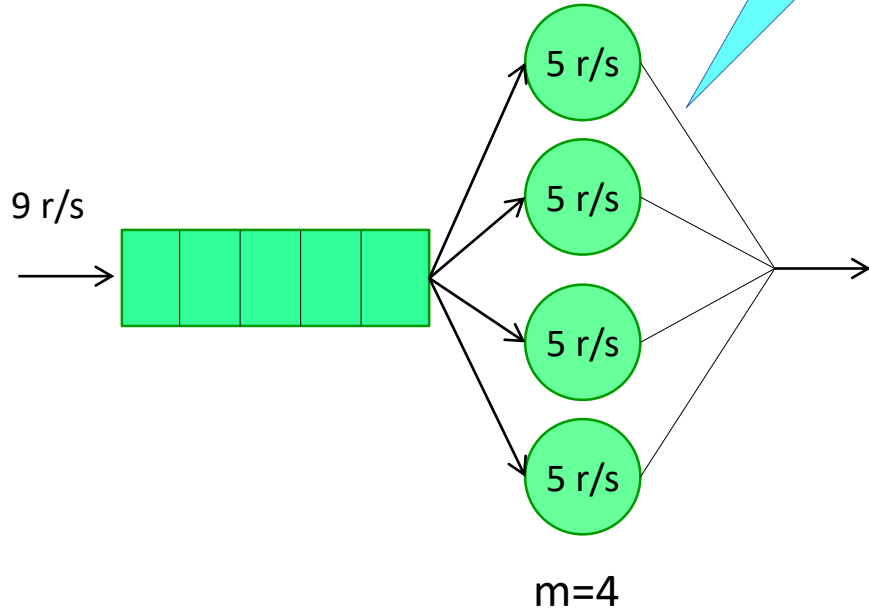
$$E[wq]$$

• Average response time:

$$E[w]$$

**In your report, you need to  
compare the output of the model  
to your measurements.  
Explain the difference!**

# Example 2: M/M/4



What changes for  $m=4$ ?

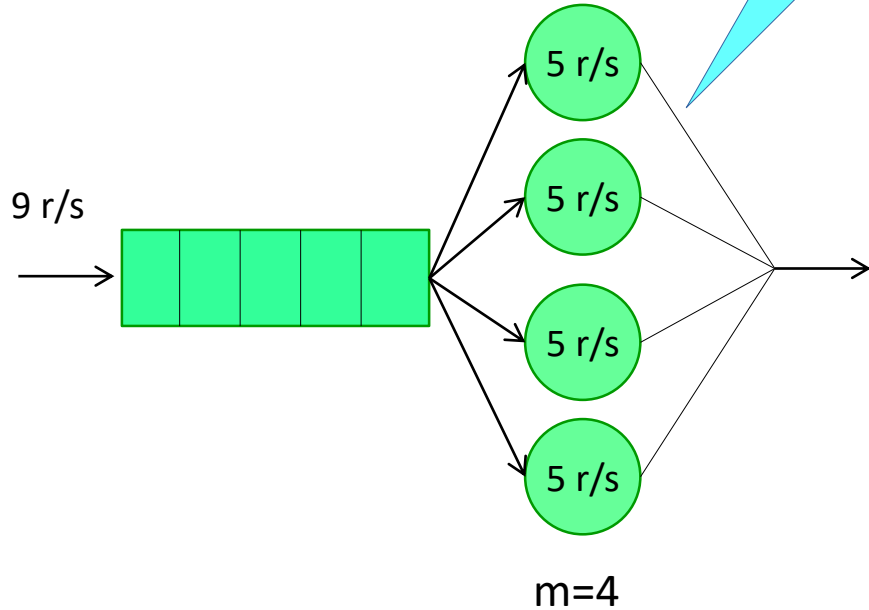
**Input:  $\lambda=9$  req/s,  $m=2$ ,  $\mu=5$  req/s**

- Utilization:
- Average service time per worker:
- Average number of requests in queue:
- Average waiting time in queue:
- Average response time:



# Example 2: M/M/4

What changes for m=4?



Input:  $\lambda=9$  req/s,  $m=2$ ,  $\mu=5$  req/s

- Utilization:  $\rho = \lambda/(m*\mu) = 45\%$
- Average service time per worker:

$$E[s] = 0.20 \text{ s}$$

- Average number of requests in queue:

$$E[nq] = 0.105$$

- Average waiting time in queue:

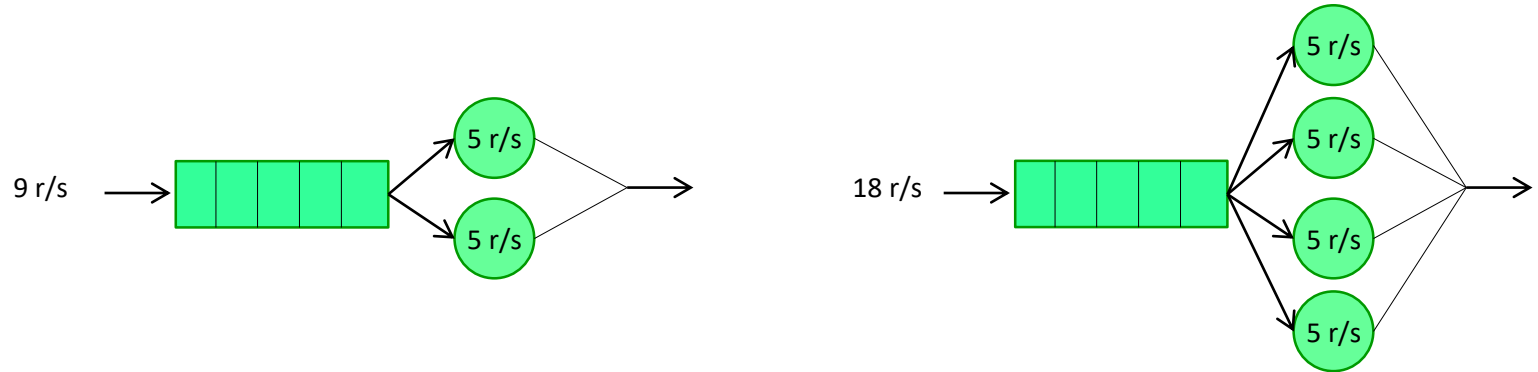
$$E[wq] = 0.012 \text{ s}$$

Lower

- Average response time:

$$E[w] = 0.012 + 0.20 = 0.212 \text{ s}$$

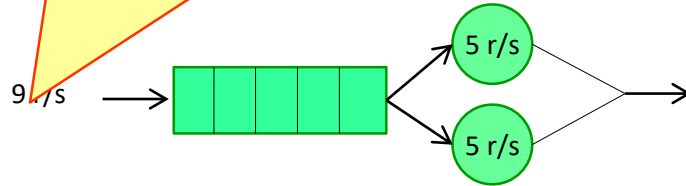
# Example 3: How do these systems differ?



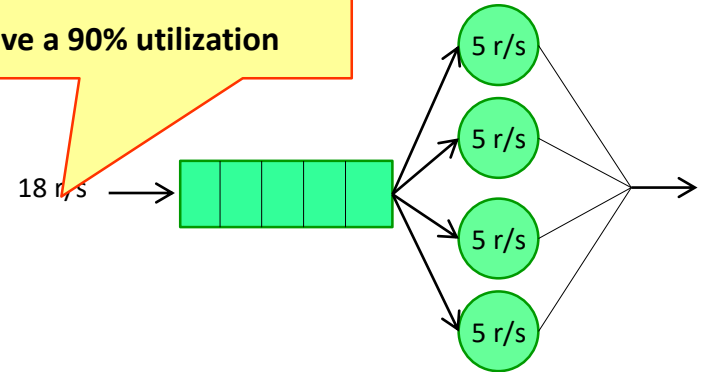
	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization		
Service time		
Queue Length		
Queuing Time		
Response Time		

# Example 3: How do these systems differ?

Both systems have a 90% utilization

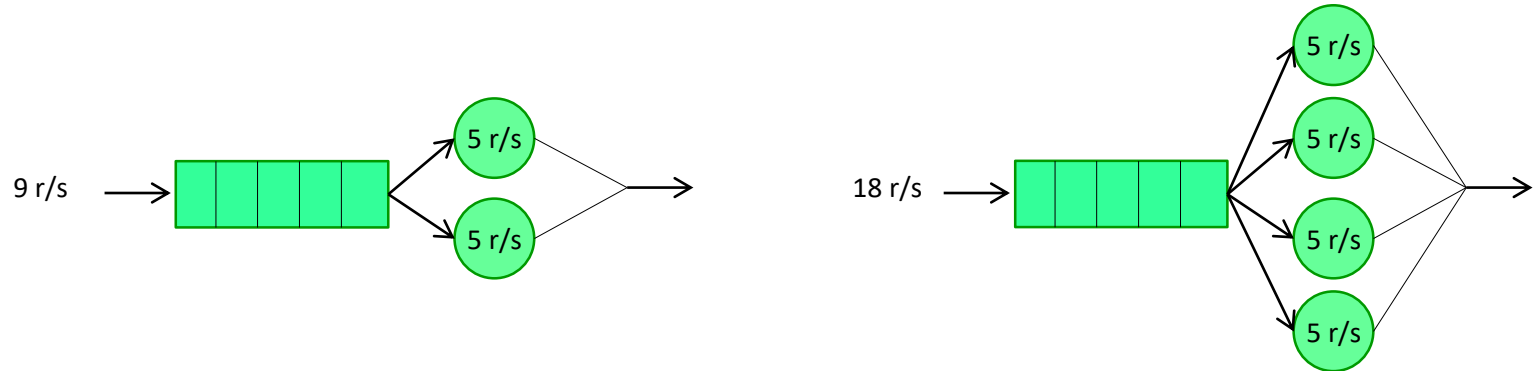


Both systems have a 90% utilization



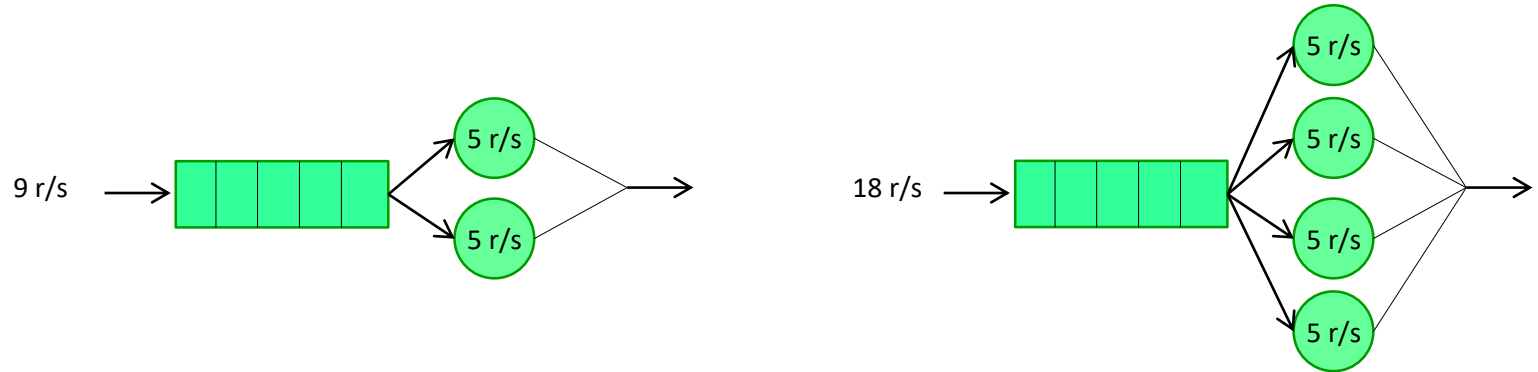
	$\lambda=9 \text{ req/s}, m=2, \mu=5 \text{ req/s}$	$\lambda=18 \text{ req/s}, m=4, \mu=5 \text{ req/s}$
Utilization	90%	90%
Service time		
Queue Length		
Queuing Time		
Response Time		

# Example 3: How do these systems differ?



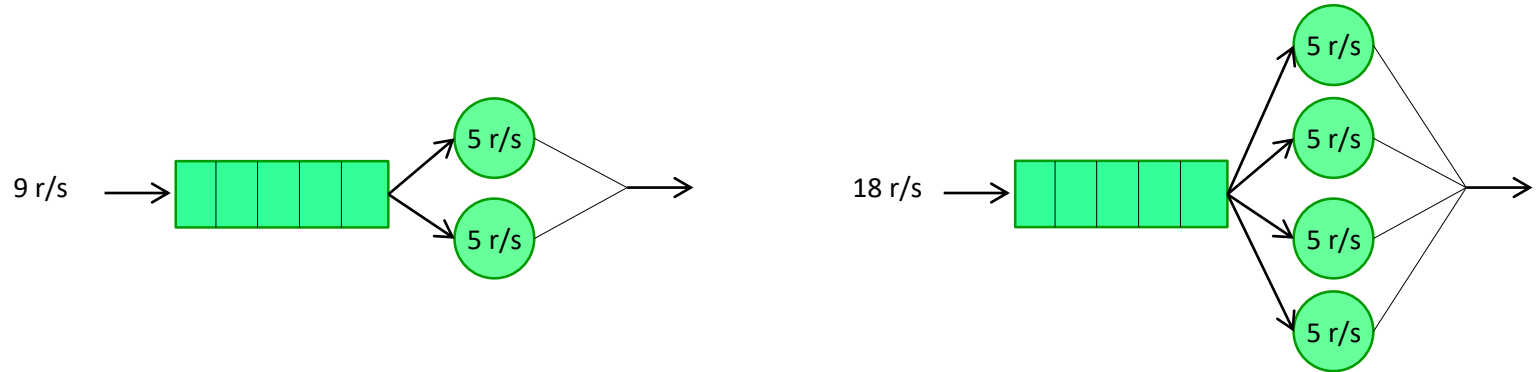
	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization	90%	90%
Service time	0.200 s	0.200 s
Queue Length		
Queuing Time		
Response Time		

# Example 3: How do these systems differ?



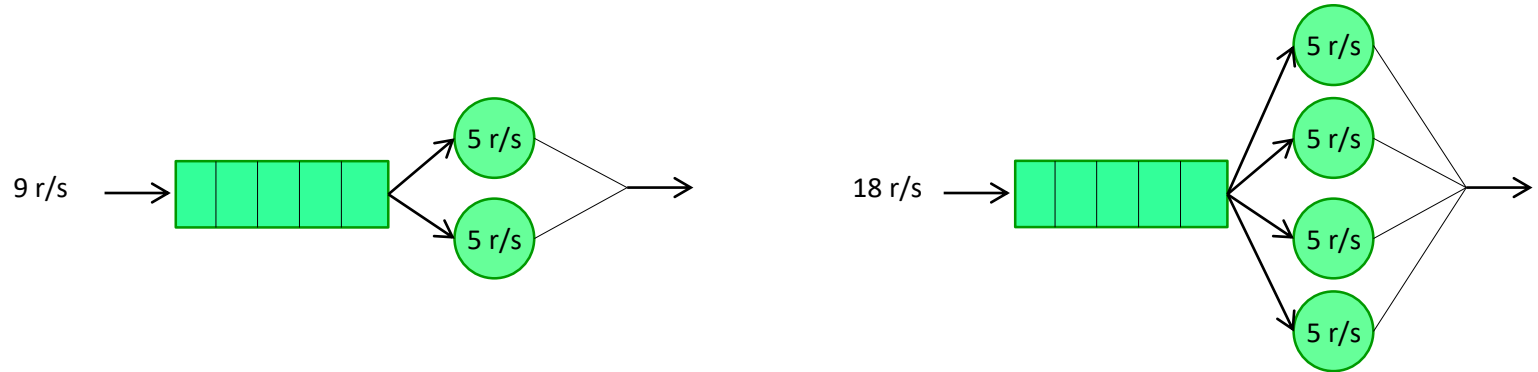
	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization	90%	90%
Service time	0.200 s	0.200 s
Queue Length	7.67 requests	7.10 requests
Queuing Time		
Response Time		

# Example 3: How do these systems differ?



	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization	90%	90%
Service time	0.200 s	0.200 s
Queue Length	7.67 requests	7.10 requests
Queuing Time	0.853 s	0.394 s
Response Time		

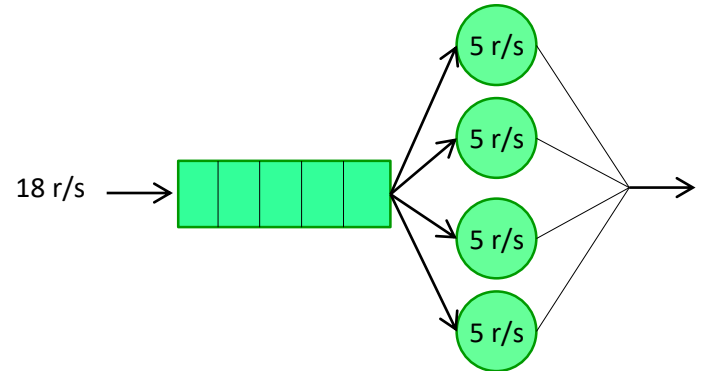
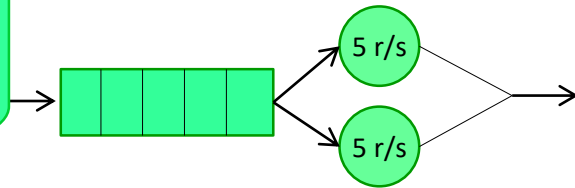
# Example 3: How do these systems differ?



	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization	90%	90%
Service time	0.200 s	0.200 s
Queue Length	7.67 requests	7.10 requests
Queuing Time	0.853 s	0.394 s
Response Time	1.053 s	0.594 s

# Example 3: How do these systems differ?

What would be better?  
2x M/M/2 or 1x M/M/4

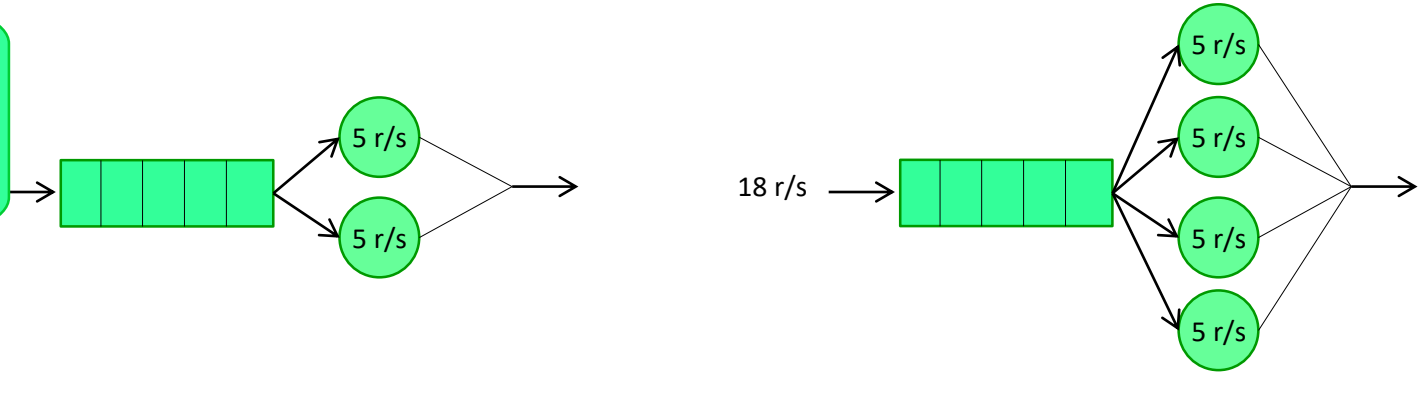


	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization	90%	90%
Service time	0.200 s	0.200 s
Queue Length	7.67 requests	7.10 requests
Queuing Time	0.853 s	0.394 s
Response Time	1.053 s	0.594 s



# Example 3: How do these systems differ?

What would be better?  
 2x M/M/2 or 1x M/M/4  
 M/M/4 has less queueing!



	$\lambda=9$ req/s, $m=2$ , $\mu=5$ req/s	$\lambda=18$ req/s, $m=4$ , $\mu=5$ req/s
Utilization	90%	90%
Service time	0.200 s	0.200 s
Queue Length	7.67 requests	7.10 requests
Queuing Time	0.853 s	0.394 s
Response Time	1.053 s	0.594 s