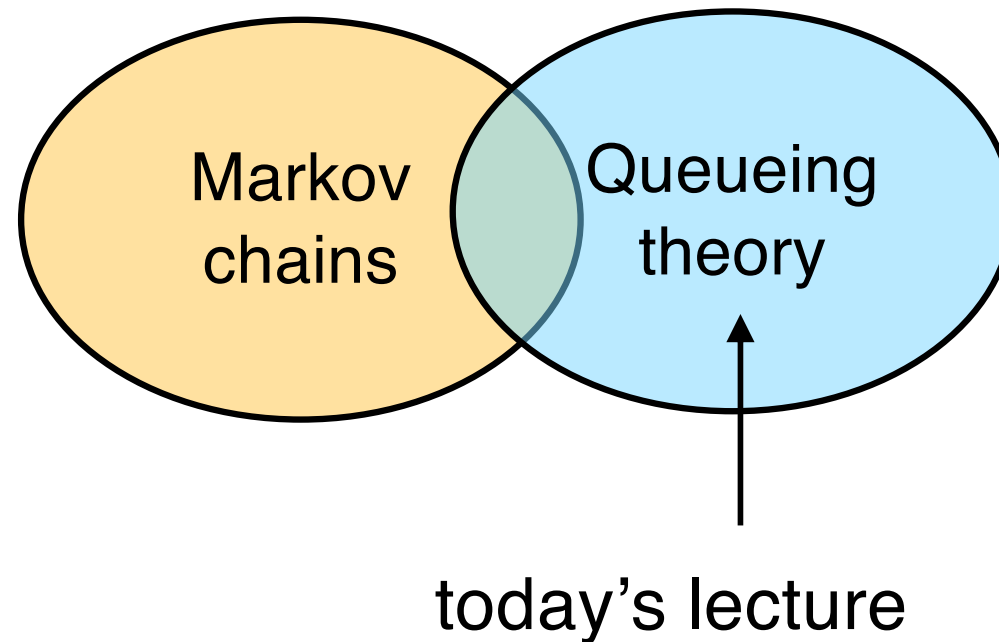


# **Performance Evaluation of Networks**

Sara Alouf

# Ch 5 – The General Service Time Queue



- $M / G / 1$  FIFO queue
- $M / G / 1$  FIFO queue with vacations

# $M / G / 1$ FIFO Queue

- Arrivals **Poisson** process rate  $\lambda$
- Service time  $\rightarrow$  **General** distribution and independence
  - ▶ Service times are independent identically distributed
  - ▶  $\sigma$  generic service time

$$G(x) = P(\sigma \leq x), \quad x \geq 0$$

$$E[\sigma] = \int_0^{\infty} (1 - G(x)) dx = \frac{1}{\mu}$$

$$E[\sigma^2] = \int_0^{\infty} x^2 dG(x)$$

$$\text{Var}[\sigma] = E[\sigma^2] - (E[\sigma])^2$$

- **First-in-first-out** service discipline

# *M / G / 1* FIFO Queue

- Load  $\rho = \frac{\lambda}{\mu}$
- Queue size  $N(t) \rightarrow$  not Markov chain  
sojourn time in a state is not Exp()
- Expected waiting time in steady-state  $\bar{W}$
- Pollaczek-Khinchin formula

$$\bar{W} = \frac{\lambda E[\sigma^2]}{2(1 - \rho)} = \frac{\rho}{2(1 - \rho)} \cdot \frac{\text{Var}(\sigma) + E[\sigma]^2}{E[\sigma]}$$

- Higher service time variability  $\rightarrow$  longer waiting times

# Proof of Pollaczek-Khinchin Formula

- For customer  $i$

- ▶ Arrival time  $t_i$
- ▶ Service time  $\sigma_i$
- ▶ Waiting time  $W_i$

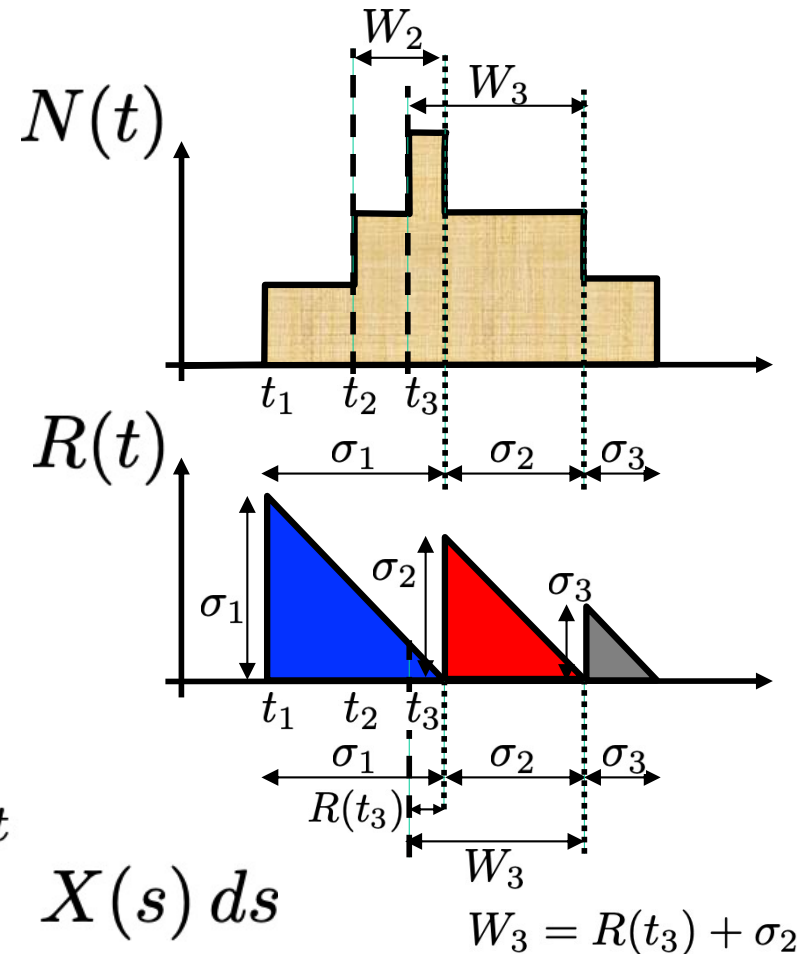
- Number of customers waiting

- ▶ at time  $t \rightarrow X(t)$
- ▶ at time  $t_i \rightarrow X(t_i) = X(t_i^-)$

- ▶ expectation  $\bar{X} = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t X(s) ds$

- Residual service time

- ▶ at time  $t \rightarrow R(t)$
- ▶ expectation  $\bar{R} = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t R(s) ds$



# Proof of Pollaczek-Khinchin Formula

- Customer  $i$  sees  $X(t_i)$  customers waiting

$$W_i = R(t_i) + \sigma_{i-1} + \sigma_{i-2} + \dots + \sigma_{i-X(t_i)}$$

$$E[W_i] = E[R(t_i)] + E \left[ \sum_{j=1}^{X(t_i)} \sigma_{i-j} \right]$$

- $X(t_i)$  consists of customers  $i-1, \dots, i-X(t_i)$   
they have not been served yet  
→  $X(t_i)$  independent of their service time
- Use Wald's formula

$$\underbrace{E[W_i]}_{\overline{W}} = \underbrace{E[R(t_i)]}_{??} + \underbrace{E[X(t_i)]}_{??} E[\sigma]$$

# Proof of Pollaczek-Khinchin Formula

- Take limit  $i \rightarrow \infty$  and use PASTA property
- We have

$$\lim_{i \rightarrow \infty} E[R(t_i)] = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t R(s) ds = \overline{R}$$

expected residual service time at **arrival epochs** in steady state

= **time average** of residual service time

- Similarly

$$\lim_{i \rightarrow \infty} E[X(t_i)] = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t X(s) ds = \overline{X}$$

- Therefore

$$E[W_i] = E[R(t_i)] + E[X(t_i)]E[\sigma] \rightarrow \overline{W} = \overline{R} + \frac{\overline{X}}{\mu}$$

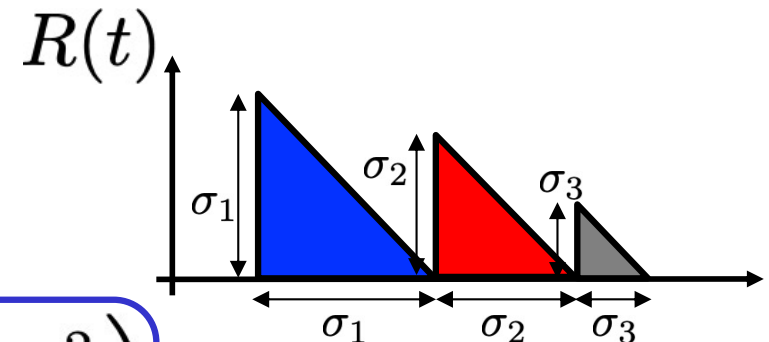
# Proof of Pollaczek-Khinchin Formula

- If  $\rho < 1$  (stability condition of  $M / G / 1$ )  
 → queue empties infinitely often
- Let 0 and  $C$  be two times when system is empty
- Let  $k$  be number of customers served in  $(0, C)$
- Expected residual service time

$$\bar{R} = \lim_{C \rightarrow \infty} \frac{1}{C} \sum_{i=1}^k \frac{\sigma_i^2}{2}$$

$$= \lim_{\substack{C \rightarrow \infty \\ k \rightarrow \infty}} \left( \frac{k}{C} \right) \lim_{\substack{C \rightarrow \infty \\ k \rightarrow \infty}} \left( \frac{1}{k} \sum_{i=1}^k \frac{\sigma_i^2}{2} \right)$$

$$= \lambda \frac{E[\sigma^2]}{2}$$





# Proof of Pollaczek-Khinchin Formula

- Apply Little's formula on waiting room

$$\bar{X} = \lambda \bar{W}$$

- Recall  $\bar{W} = \bar{R} + \frac{\bar{X}}{\mu} = \bar{R} + \frac{\lambda}{\mu} \bar{W}$

$$\Rightarrow \bar{W}(1 - \rho) = \bar{R}$$

$$\Leftrightarrow \bar{W} = \frac{\bar{R}}{1 - \rho}$$

$$\Leftrightarrow \bar{W} = \frac{\lambda E[\sigma^2]}{2(1 - \rho)}$$

# *M / G / 1* FIFO Queue

- Expected waiting time

$$\overline{W} = \frac{\lambda E[\sigma^2]}{2(1 - \rho)}$$

- Expected sojourn time

$$\overline{T} = \overline{W} + \frac{1}{\mu} = \frac{1}{\mu} + \frac{\lambda E[\sigma^2]}{2(1 - \rho)}$$

- Expected number of customers waiting

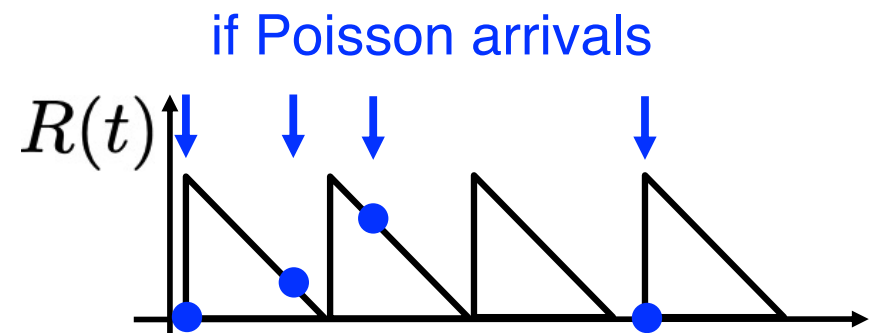
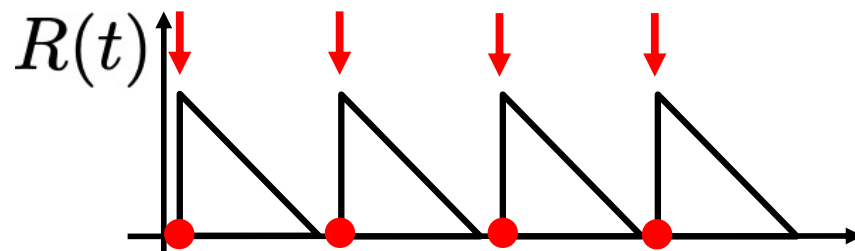
$$\overline{X} = \lambda \overline{W} = \frac{\lambda^2 E[\sigma^2]}{2(1 - \rho)}$$

- Expected queue size

$$\overline{N} = \overline{X} + \rho = \lambda \overline{T} = \rho + \frac{\lambda^2 E[\sigma^2]}{2(1 - \rho)}$$

# Example When PASTA Not True

- Consider  $D / D / 1$  FIFO queue
- Arrivals every second  $\rightarrow \lambda = 1 \text{ s}^{-1}$
- Service time 0.9 second  $\rightarrow \mu = 1/0.9 \text{ s}^{-1}$
- Load is very high  $\rightarrow \rho = 0.9$



- Time average

$$\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t R(s) ds = \lambda \frac{E[\sigma^2]}{2} = \frac{(0.9)^2}{2} = 0.405$$

- Average at arrival epochs  $\lim_{i \rightarrow \infty} E[R(t_i)] = 0$

# *M / M / 1* Versus *M / D / 1* FIFO

- Consider two queues
  - ▶ Poisson arrival rate  $\lambda$
  - ▶ One server
  - ▶ Infinite waiting room
  - ▶ FIFO service discipline
- Different service time distribution but same average  $1/\mu$ 
  - ▶ *M / M / 1* queue :  $\text{Exp}(\mu)$   $E[\sigma^2] = \frac{2}{\mu^2}$
  - ▶ *M / D / 1* queue :  $\sigma = 1/\mu$   $E[\sigma^2] = \frac{1}{\mu^2}$
- Expected waiting time

$$\bar{W}_{M/M/1} = \frac{\lambda E[\sigma^2]}{2(1-\rho)} = \frac{\lambda}{\mu^2(1-\rho)} \quad \boxed{\bar{W}_{M/D/1} = \frac{\bar{W}_{M/M/1}}{2}}$$

# *M / G / 1* FIFO Queue With Vacations

- *M / G / 1* FIFO but if queue empty: server → vacation
  - ▶ Maintenance, background task, sleep mode, power off
- First-in-first-out service discipline
- Arrivals Poisson process rate  $\lambda$
- Service time → General distribution and independence
  - ▶ Service times are independent identically distributed
  - ▶  $\sigma$  generic service time

$$G(x) = P(\sigma \leq x), \quad x \geq 0$$

$$E[\sigma] = \int_0^{\infty} (1 - G(x)) dx = \frac{1}{\mu}$$

$$E[\sigma^2] = \int_0^{\infty} x^2 dG(x)$$

# $M / G / 1$ FIFO Queue With Vacations

- Vacation time → General distribution and independence
  - ▶ vacation durations are independent identically distributed
  - ▶  $V$  generic vacation duration

$$F(x) = P(V \leq x), \quad x \geq 0$$

$$E[V] = \int_0^{\infty} (1 - F(x)) dx$$

$$E[V^2] = \int_0^{\infty} x^2 dF(x)$$

- Load  $\rho = \frac{\lambda}{\mu}$

- Queue size not Markov chain (sojourn time not Exp() )

# $M / G / 1$ FIFO Queue With Vacations

- If  $\rho < 1$  (stability condition) expected waiting time steady-state

$$\begin{aligned}\bar{W} &= \frac{\lambda E[\sigma^2]}{2(1-\rho)} + \frac{E[V^2]}{2E[V]} \\ &= \bar{W}_{M/G/1} + \frac{E[V^2]}{2E[V]}\end{aligned}$$

- Higher vacations variability  $\rightarrow$  longer waiting times

$$\frac{E[V^2]}{2E[V]} = \frac{\text{Var}[V]}{2E[V]} + \frac{E[V]}{2}$$

- To lessen impact of vacations  $\rightarrow V$  deterministic, small
- If cost to go on vacation  $\rightarrow$  tradeoff to be found

15 minutes break



# $M / G / 1$ FIFO Queue With Vacations

- Expected waiting time in steady-state  $\overline{W}$

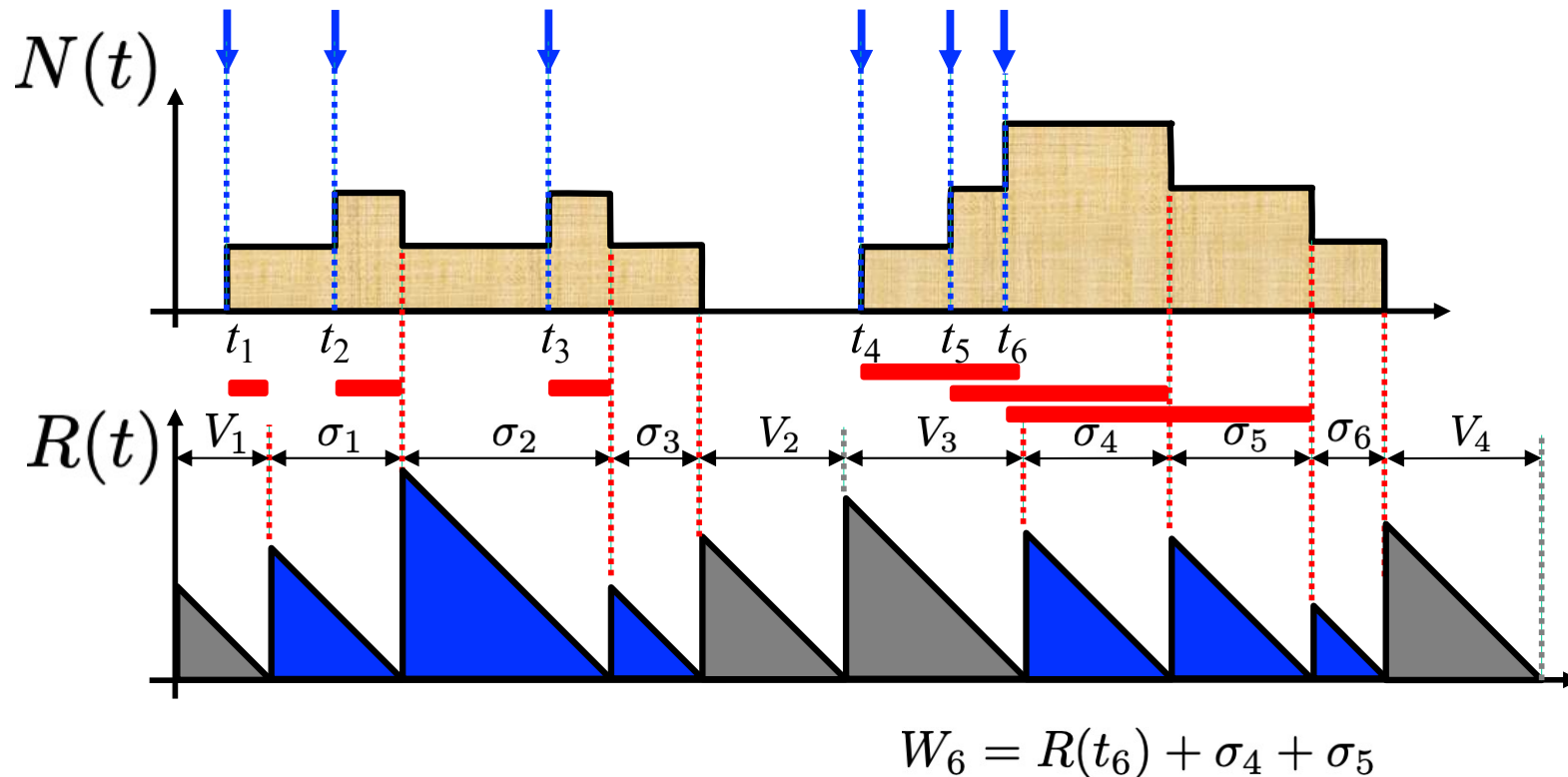
$$\overline{W} = \frac{\lambda E[\sigma^2]}{2(1 - \rho)} + \frac{E[V^2]}{2E[V]}$$

# Proof

- For customer  $i$ 
  - ▶ Arrival time  $t_i$
  - ▶ Service time  $\sigma_i$
  - ▶ Waiting time  $W_i$
- Number of customers waiting
  - ▶ at time  $t \rightarrow X(t)$
  - ▶ at time  $t_i \rightarrow X(t_i) = X(t_i^-)$
  - ▶ expectation  $\bar{X} = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t X(s) ds$
- $k$ th server vacation time  $V_k$

# Proof

- Residual time at server at time  $t \rightarrow R(t)$ 
  - ▶ if server **busy**  $\rightarrow$  residual service time
  - ▶ if server in **vacation**  $\rightarrow$  residual vacation time



# Proof

- Customer  $i$  sees  $X(t_i)$  customers waiting

$$W_i = R(t_i) + \sigma_{i-1} + \sigma_{i-2} + \dots + \sigma_{i-X(t_i)}$$

$$E[W_i] = E[R(t_i)] + E \left[ \sum_{j=1}^{X(t_i)} \sigma_{i-j} \right]$$

- Use Wald's formula ( $X(t_i)$  independent of all  $\sigma_{i-j}$ )

$$E[W_i] = E[R(t_i)] + E[X(t_i)]E[\sigma]$$

- Take limit  $i \rightarrow \infty$  and use PASTA property

$$\left. \begin{aligned} \overline{W} &= \overline{R} + \frac{\overline{X}}{\mu} \end{aligned} \right\} \Rightarrow$$

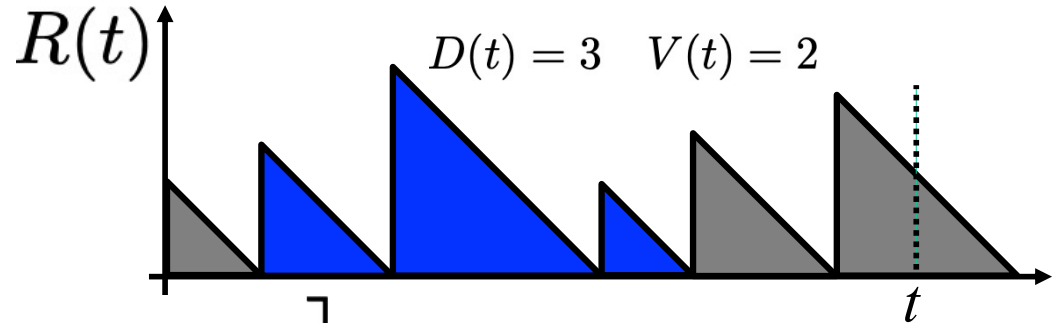
$$\overline{W} = \frac{\overline{R}}{1 - \rho}$$

same  
expression  
different  
 $\overline{R}$

- By Little's formula  $\overline{X} = \lambda \overline{W}$

# Proof: Find $\bar{R}$

- $\rho < 1 \rightarrow$  queue will empty infinitely often
- $D(t) \rightarrow$  number of customers **fully** served in  $(0, t)$
- $V(t) \rightarrow$  number of **complete** vacations in  $(0, t)$
- Expected residual time



$$\bar{R} = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t R(u) du$$

$$= \lim_{t \rightarrow \infty} \frac{1}{t} \left[ \sum_{i=1}^{D(t)} \frac{\sigma_i^2}{2} + \sum_{k=1}^{V(t)} \frac{V_k^2}{2} + \text{trapezoid} \right]$$

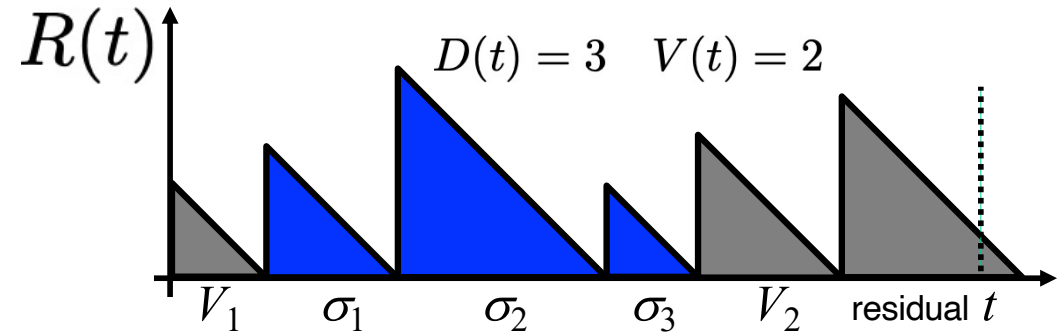
$$= \lim_{t \rightarrow \infty} \left[ \frac{D(t)}{t} \frac{1}{D(t)} \sum_{i=1}^{D(t)} \frac{\sigma_i^2}{2} + \frac{V(t)}{t} \frac{1}{V(t)} \sum_{k=1}^{V(t)} \frac{V_k^2}{2} + \frac{\text{trapezoid}}{t} \right]$$

$$= \lambda \frac{E[\sigma^2]}{2} + \left( \lim_{t \rightarrow \infty} \frac{V(t)}{t} \right) \frac{E[V^2]}{2}$$

# Proof: Find $\lim_{t \rightarrow \infty} \frac{V(t)}{t}$

■ We have

$$t = \sum_{i=1}^{D(t)} \sigma_i + \sum_{k=1}^{V(t)} V_k + \text{residual}$$



$$1 = \frac{D(t)}{t} \frac{1}{D(t)} \sum_{i=1}^{D(t)} \sigma_i + \frac{V(t)}{t} \frac{1}{V(t)} \sum_{k=1}^{V(t)} V_k + \frac{\text{residual}}{t}$$

take  
limit  
 $t \rightarrow \infty$

$$1 = \lambda E[\sigma] + \left( \lim_{t \rightarrow \infty} \frac{V(t)}{t} \right) E[V]$$

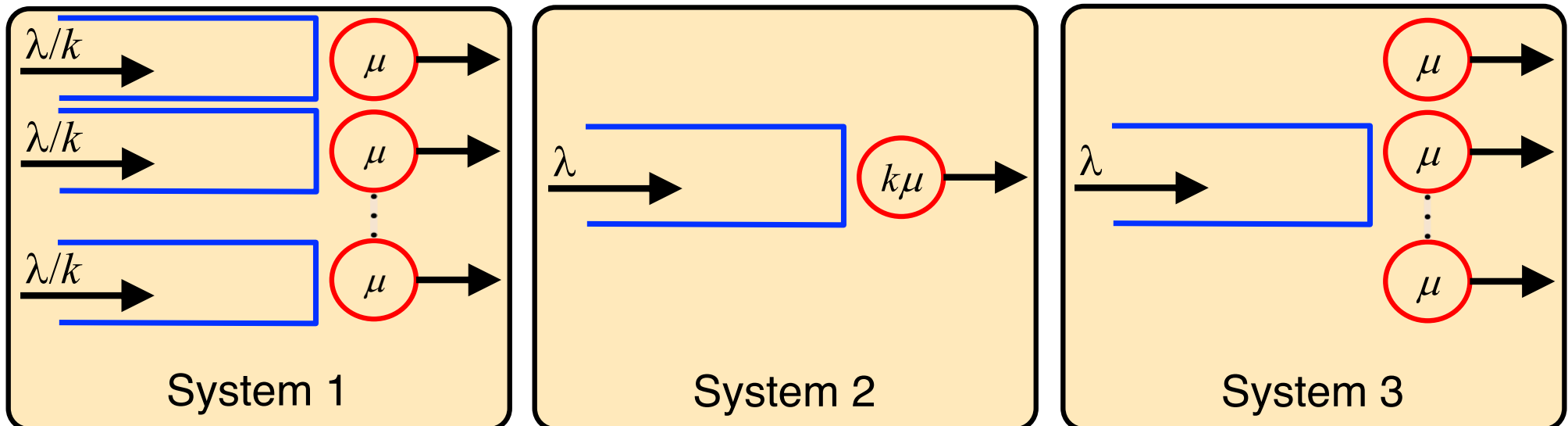
$$\Rightarrow \lim_{t \rightarrow \infty} \frac{V(t)}{t} = \frac{1 - \rho}{E[V]}$$

# Proof: Recap

- $\lim_{t \rightarrow \infty} \frac{V(t)}{t} = \frac{1 - \rho}{E[V]}$
- $$\begin{aligned}\bar{R} &= \lambda \frac{E[\sigma^2]}{2} + \left( \lim_{t \rightarrow \infty} \frac{V(t)}{t} \right) \frac{E[V^2]}{2} \\ &= \lambda \frac{E[\sigma^2]}{2} + \left( \frac{1 - \rho}{E[V]} \right) \frac{E[V^2]}{2}\end{aligned}$$
- $$\begin{aligned}\bar{W} &= \frac{\bar{R}}{1 - \rho} \\ &= \frac{\lambda E[\sigma^2]}{2(1 - \rho)} + \frac{E[V^2]}{2E[V]} \quad \checkmark \\ &= \bar{W}_{M/G/1} + \text{effect of vacations}\end{aligned}$$

# Exercise: Compare Different Organizations

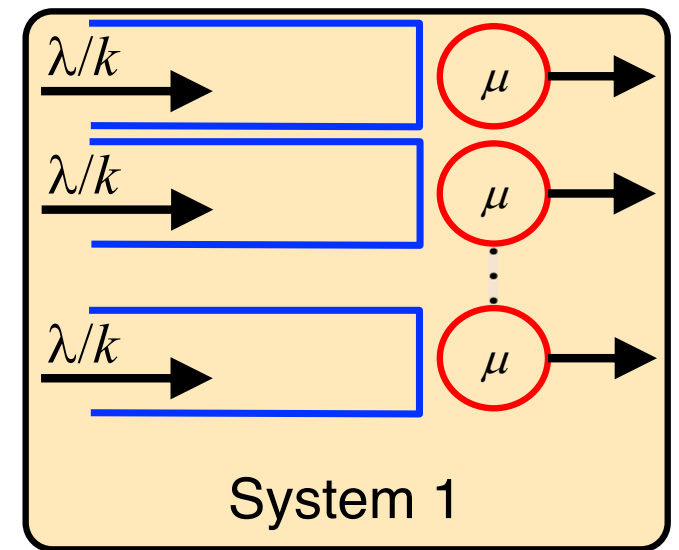
- Global traffic Poisson rate  $\lambda$
  - Total service rate  $k\mu$
  - Service time is exponentially distributed
  - Objective: compare three systems organizations
- $\left. \begin{array}{l} \text{Global traffic Poisson rate } \lambda \\ \text{Total service rate } k\mu \end{array} \right\} \rho = \frac{\lambda}{k\mu}$



- Order systems according to expected sojourn time



# System 1



- Each queue is  $M / M / 1$  queue
- Arrival rate  $\lambda/k$
- Service rate  $\mu$
- Infinite queue: stability condition  $\lambda/k < \mu$  ( $\rho < 1$ )

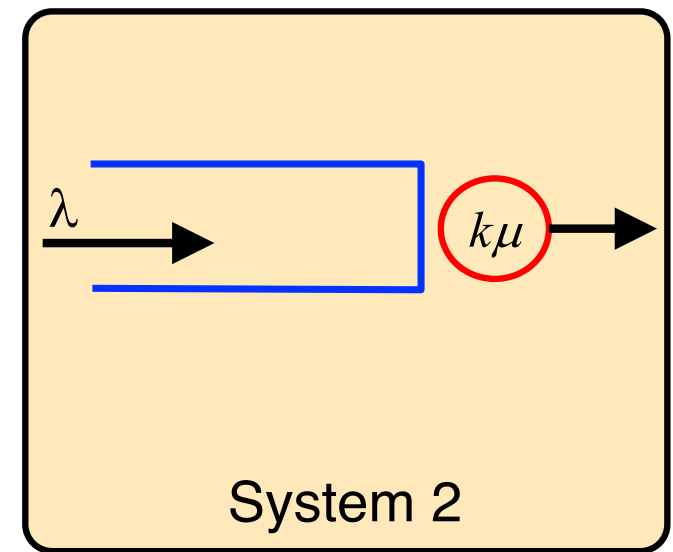
- Expected queue size in one queue  $\frac{\lambda/k}{\mu - \lambda/k}$

- Expected number of customers in System 1

$$\bar{N}_1 = \frac{\lambda}{\mu - \lambda/k} = \frac{k\rho}{1 - \rho}$$

- By Little's formula  $\bar{T}_1 = \frac{\bar{N}_1}{\lambda} = \frac{1}{\mu - \lambda/k} = \frac{k}{k\mu - \lambda}$

## System 2



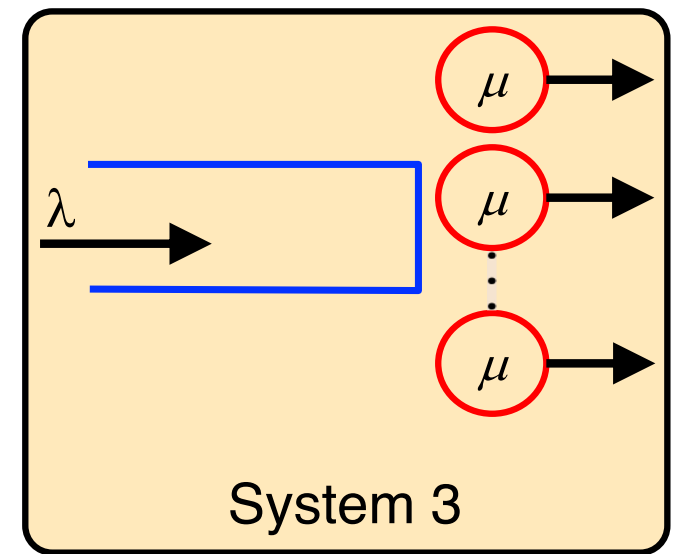
- Queue is  $M / M / 1$  queue
- Arrival rate  $\lambda$
- Service rate  $k\mu$
- Infinite queue: stability condition  $\lambda < k\mu$  ( $\rho < 1$ )
- Expected number of customers in System 2

$$\bar{N}_2 = \frac{\lambda}{k\mu - \lambda} = \frac{\rho}{1 - \rho}$$

- By Little's formula

$$\bar{T}_2 = \frac{\bar{N}_2}{\lambda} = \frac{1}{k\mu - \lambda}$$

# System 3



- Queue is  $M / M / k$  queue
- Arrival rate  $\lambda$
- Service rate  $\mu$
- Infinite queue: stability condition  $\lambda < k\mu$  ( $\rho < 1$ )
- Expected number of customers in System 3

$$\bar{N}_3 = \bar{N}_{\text{wait}} + \bar{N}_{\text{servers}}$$

- By Little's formula on all servers

$$\bar{N}_{\text{servers}} = \lambda \frac{1}{\mu}$$

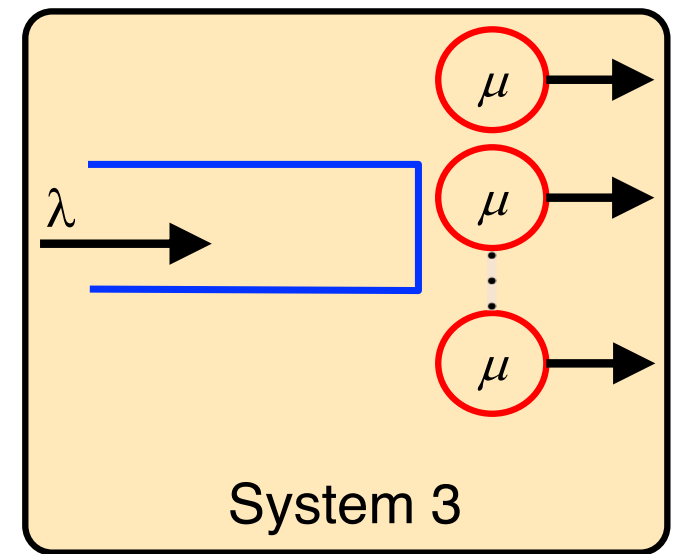
- Let  $X_{\text{wait}}$  number of customers waiting in queue

$$\bar{N}_{\text{wait}} = E[X_{\text{wait}}]$$

law of total  
probabilities

$$= E[X_{\text{wait}} | \text{wait}] P_{\text{wait}} + E[X_{\text{wait}} | \text{no wait}] (1 - P_{\text{wait}})$$

# System 3



Conditioning on fact all servers busy  
 $X_{\text{wait}}$  same as queue size in  $M / M / 1$   
 queue with arrival rate  $\lambda$  service rate  $k\mu$

$$E[X_{\text{wait}} | \text{wait}] = \bar{N}_2$$

- All servers are busy with probability (see lecture 4)

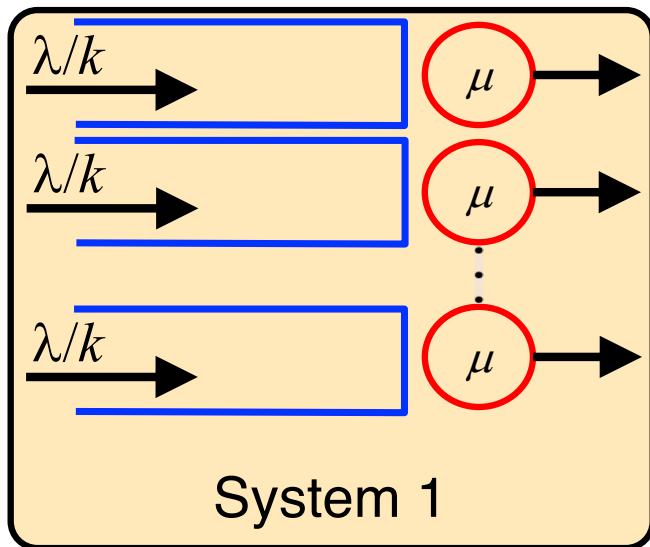
$$P_{\text{wait}} = \frac{\pi_0 (k\rho)^k}{k!(1-\rho)} = \frac{\frac{(k\rho)^k}{k!(1-\rho)}}{\sum_{i=0}^{k-1} \frac{(k\rho)^i}{i!} + \frac{(k\rho)^k}{k!(1-\rho)}}$$

- Expected number of customers in System 3

$$\bar{N}_3 = \frac{\lambda}{k\mu - \lambda} P_{\text{wait}} + \frac{\lambda}{\mu}$$

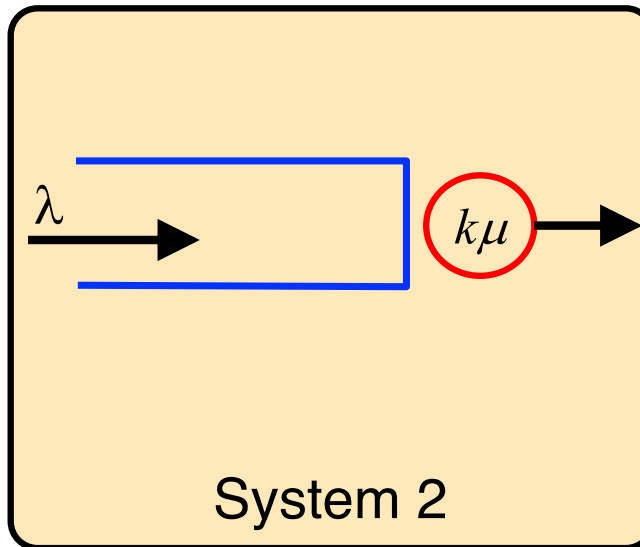
- By Little's formula  $\bar{T}_3 = \frac{1}{k\mu - \lambda} P_{\text{wait}} + \frac{1}{\mu}$

# Exercise: Compare Different Organizations



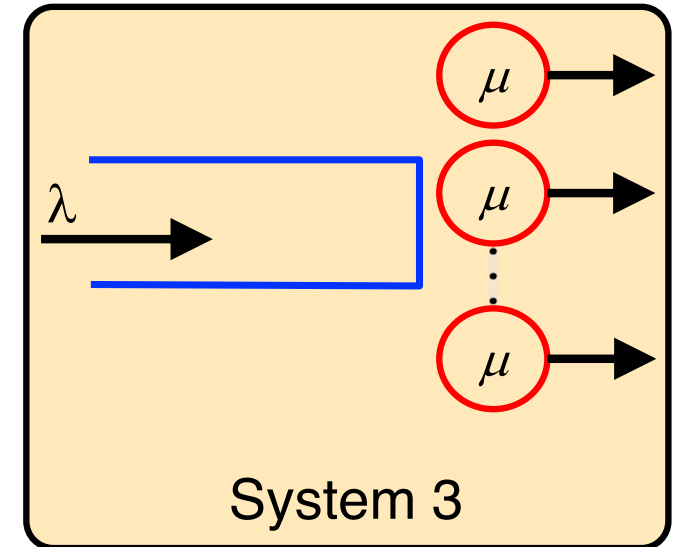
$$\bar{N}_1 = \frac{k\rho}{1-\rho}$$

$$\bar{T}_1 = \frac{k}{k\mu - \lambda}$$



$$\bar{N}_2 = \frac{\rho}{1-\rho}$$

$$\bar{T}_2 = \frac{1}{k\mu - \lambda}$$



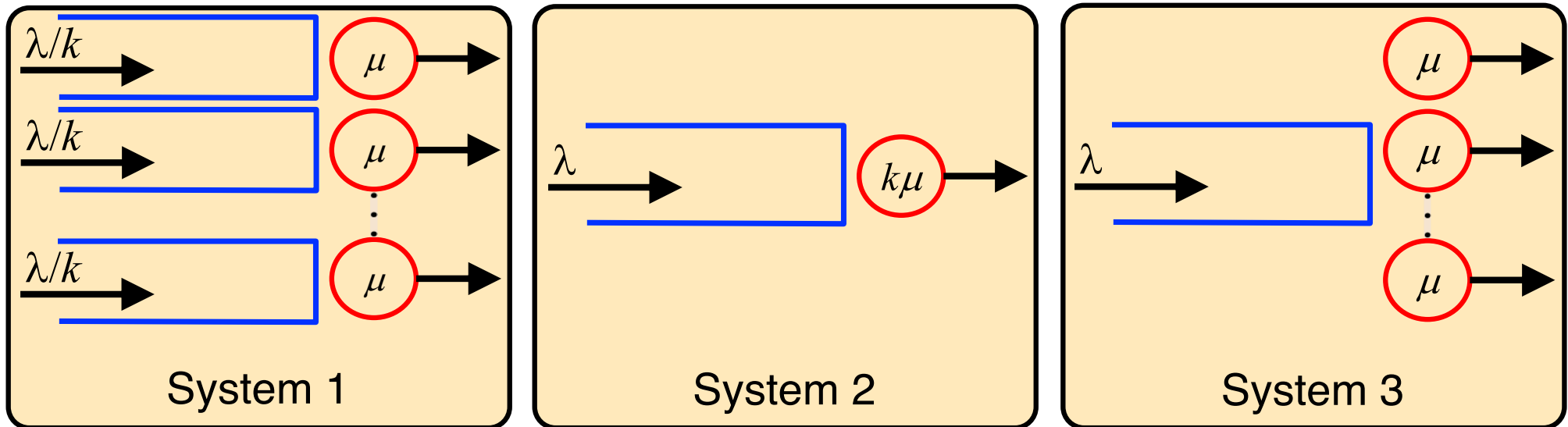
$$\bar{N}_3 = \bar{N}_2 P_{\text{wait}} + \frac{\lambda}{\mu}$$

$$\bar{T}_3 = \bar{T}_2 P_{\text{wait}} + \frac{1}{\mu}$$

- System 2 is  $k$  times better than System 1

- What about System 3?  $\frac{\bar{T}_3}{\bar{T}_2} = P_{\text{wait}} + k(1-\rho) > 1$

# Exercise: Compare Different Organizations



- What about System 3?  $\frac{\bar{T}_3}{\bar{T}_2} = P_{\text{wait}} + k(1 - \rho) > 1$
- If very low utilization  $\rightarrow$  ratio close to  $k$

System 3 **almost  $k$  times worse** than System 2

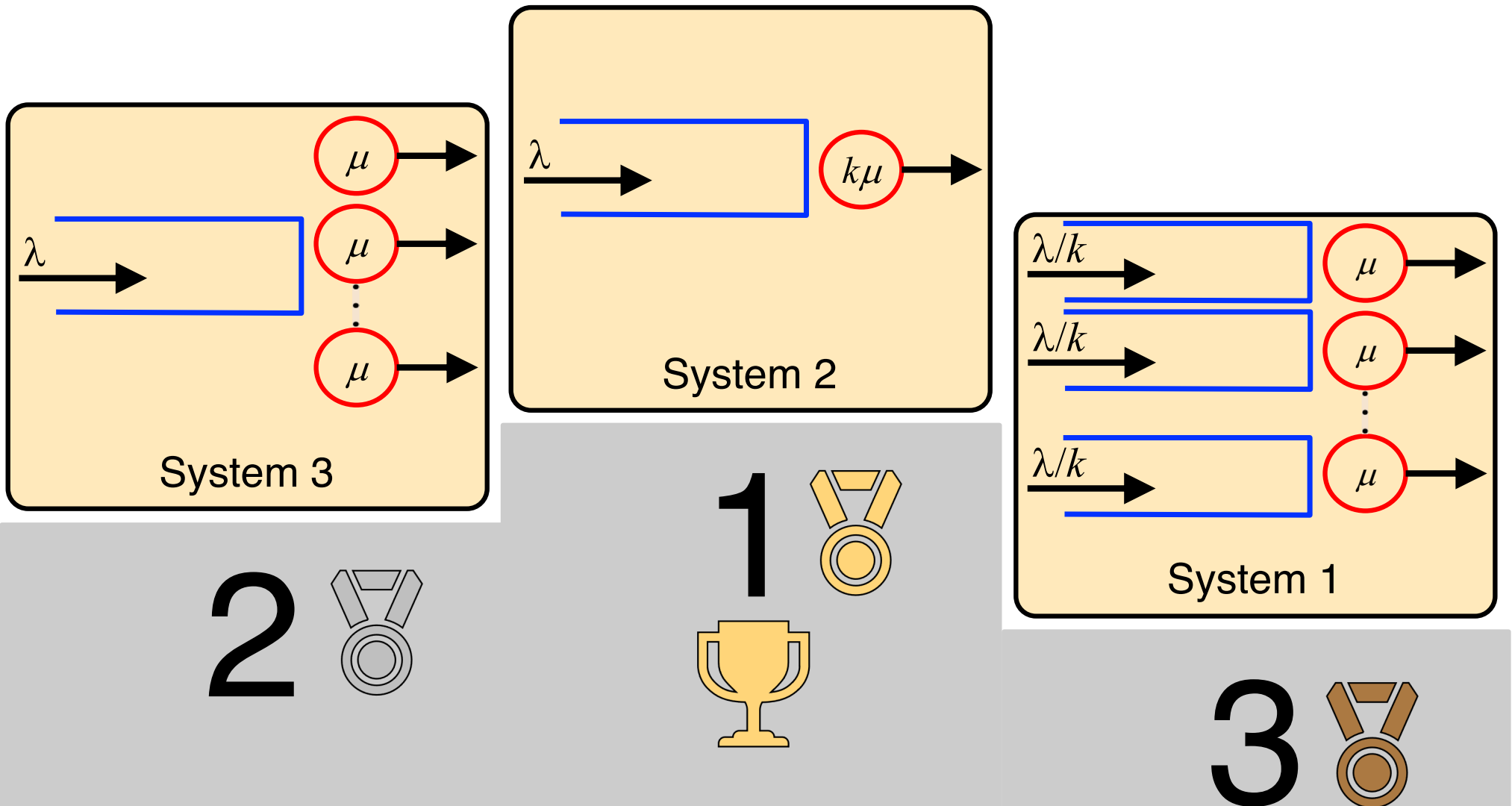
$\rightarrow$  System 3 slightly better than System 1

- If very high utilization  $\rightarrow$  ratio close to 1

System 3 **almost same** (slightly worse) as System 2

Conclusion  $\bar{T}_2 < \bar{T}_3 < \bar{T}_1$

- ... and the winner is: System 2!



# For next week

- Lesson 5 to revise
- Homework 5 to return on Tuesday 15 October before 9 am
- Lesson 6 to read before Lecture 6