#### **FINDING THE OPTIMAL QUANTUM SIZE** Revisiting the M/G/1 Round-Robin Queue

VARUN GUPTA Carnegie Mellon University







- arrival rate =  $\lambda$
- job sizes i.i.d. ~ S
- load  $\rho = \lambda E[S]$
- $C^2 = \frac{var(S)}{E[S]^2}$

Squared coefficient of variability (SCV) of job sizes:  $C^2 \ge 0$ 



#### GOAL

#### Optimal operating q for M/G/1/RR with overheads and high C<sup>2</sup>

effect of preemption overheads

#### **SUBGOAL**

Sensitivity Analysis: Effect of q and C<sup>2</sup> on M/G/1/RR performance

(no switching overheads)

#### **PRIOR WORK**

 Lots of *exact* analysis: [Wolff70], [Sakata et al.71], [Brown78]

**No closed-form solutions/bounds** 

**\Im** No simple expressions for interplay of q and  $C^2$ 

- $\Box$  Effect of q and  $C^2$  on mean response time
  - Approximate analysis
  - □ Bounds for M/G/1/RR

No preemption overheads

Choosing the optimal quantum size

Effect of preemption overheads

# Approximate sensitivity analysis of M/G/1/RR

Approximation assumption 1:

Service quantum ~ Exp(1/q)



Job size distribution 
$$\sim \begin{cases} 0 & \text{w.p. } p \\ \text{Exp}(\mu) & \text{w.p. } 1-p \end{cases}$$

$$E[T^{RR*}] = E[T^{PS}] \left[ 1 + \frac{C^2 - 1}{C^2 + 1} \cdot \frac{\lambda}{\frac{1}{q} + \frac{2}{C^2 + 1} \frac{1}{E[S]}} \right]$$

# Approximate sensitivity analysis of M/G/1/RR

$$E[T^{RR*}] = E[T^{PS}] \left[ 1 + \frac{C^2 - 1}{C^2 + 1} \cdot \frac{\lambda}{\frac{1}{q} + \frac{2}{C^2 + 1} \frac{1}{E[S]}} \right]$$

- Monotonic in q
  - Increases from  $E[T^{PS}] \rightarrow E[T^{FCFS}]$
- Monotonic in  $C^2$ 
  - Increases from  $E[T^{PS}] \rightarrow E[T^{PS}](1+\lambda q)$

#### For high *C*<sup>2</sup>: *E*[*T*<sup>RR\*</sup>] ≈ *E*[*T*<sup>PS</sup>](1+λ*q*)

- $\Box$  Effect of *q* and *C*<sup>2</sup> on mean response time
  - ✓ Approximate analysis

Bounds for M/G/1/RR

□ Choosing the optimal quantum size

## M/G/1/RR bounds

Assumption: job sizes  $\in \{0, q, \dots, Kq\}$ 



As  $K \to \infty$ : sup  $E[T^{RR}] = E[T^{PS}] \left[ 1 + \frac{(1+\rho)\lambda q}{2} \right]$ 

- ✓ Effect of q and  $C^2$  on mean response time
  - ✓ Approximate analysis
  - ✓ Bounds for M/G/1/RR
- □ Choosing the optimal quantum size

## Optimizing q

Preemption overhead = h

1. 
$$q' = q + h$$
$$E[S]' = E[S]\left(1 + \frac{h}{q}\right)$$
$$\rho' = \lambda E[S]'$$

2. Minimize  $E[T^{RR}]$  upper bound from Theorem:

$$q^* = \operatorname{argmin}_q \frac{E[S]'}{1-\rho'} \left[ 1 + \frac{(1+\rho')\lambda q'}{2} \right]$$

**Common case:** 
$$\frac{h}{E[S]} \ll (1 - \rho)$$
  
 $q^* \approx \alpha(\rho) \sqrt{hE[S]}$ 













- ✓ Effect of q and  $C^2$  on mean response time
  - ✓ Approximate analysis
  - ✓ Bounds for M/G/1/RR
- ✓ Choosing the optimal quantum size

### **Conclusion/Contributions**

- Simple approximation and bounds for M/G/1/RR
- Optimal quantum size for handling highly variable job sizes under preemption overheads



## Bounds – Proof outline

 $\boldsymbol{D}^{T} = \boldsymbol{A}_{P}\boldsymbol{D}^{T} + \boldsymbol{b}$ 

- $D_i$  = mean delay for *I*<sup>th</sup> quantum of service
- $\boldsymbol{D} = [D_1 \ D_2 \ \dots \ D_K]$
- **D** is the fixed point of a monotone linear system:

 $f_1 = 0$  $f_2 = 0$  $D_2$ Sufficient condition for upper bound: Sufficient condition D' for lower bound:  $D^{*T} \geq A_P D^{*T} + b$ **D**\*  $D^{\prime T} \leq A_P D^{\prime T} + b$ for all P D for all P 22  $D_1$ 

## Optimizing q

- Preemption overhead = h
- q' = q+h,  $E[S]' \rightarrow E[S] (1+h/q)$ ,  $\rho' = \lambda E[S]'$
- $\min_{q} \frac{E[S]'}{1-\rho'} \left[ 1 + \frac{(1+\rho')\lambda q'}{2} \right]$

Heavy traffic: 
$$\frac{1}{1-\rho} \gg \frac{E[S]}{h}$$
  
 $q^* \approx \frac{2h}{1-\rho}$ 

Small overhead:  $\frac{1}{1-\rho} \ll \frac{E[S]}{h}$  $q^* \approx K(\rho) \sqrt{hE[S]}$