**Aalto University**
**School of Electrical**
**Engineering**

# Optimal scheduling problem for scalable queues
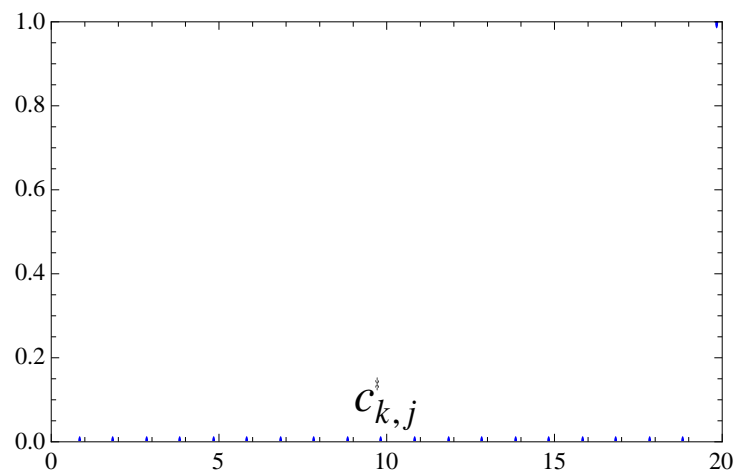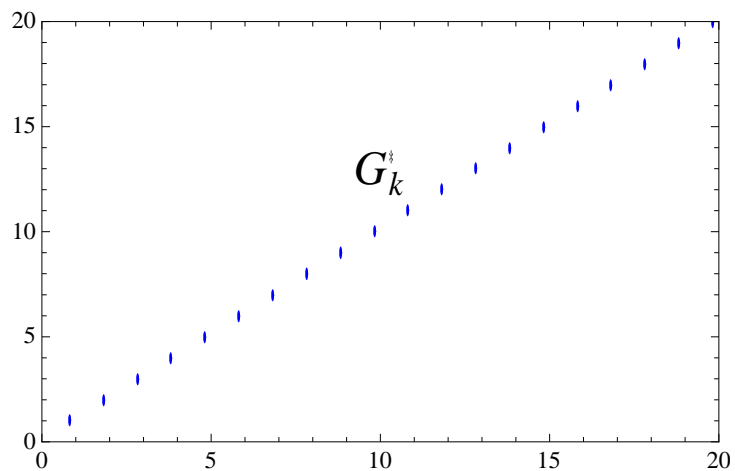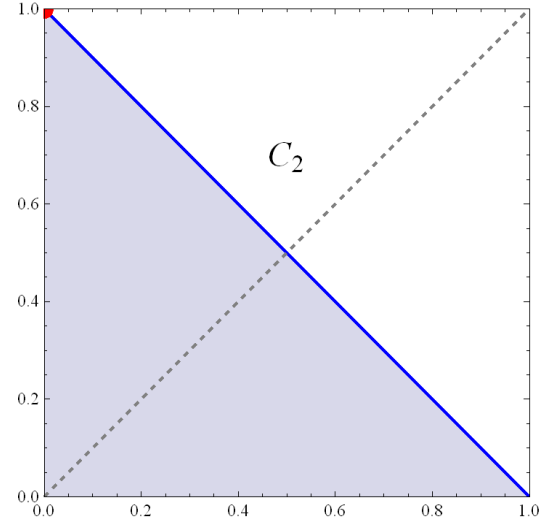
Samuli Aalto

Aalto University, Finland

# Alpha = 1.0
# (single-server queue)

# Alpha = 1.2



$C_2$



$G_k^*$



$c_{k,j}^*$

# Alpha = 2.0



$C_2$



$G_k^*$



$c_{k,j}^*$

# Alpha = 5.0

# Alpha = infinite (infinite-server queue)



$$C_2$$



$$G_k^{\ast}$$



$$c_{k,j}^{\ast}$$

# Scalable queue



$C_2$

- Service system where the service capacity scales with number of jobs

- Policy: When there are $k$ jobs with sizes

$$s_1 \geq \ldots \geq s_k$$

choose a rate vector

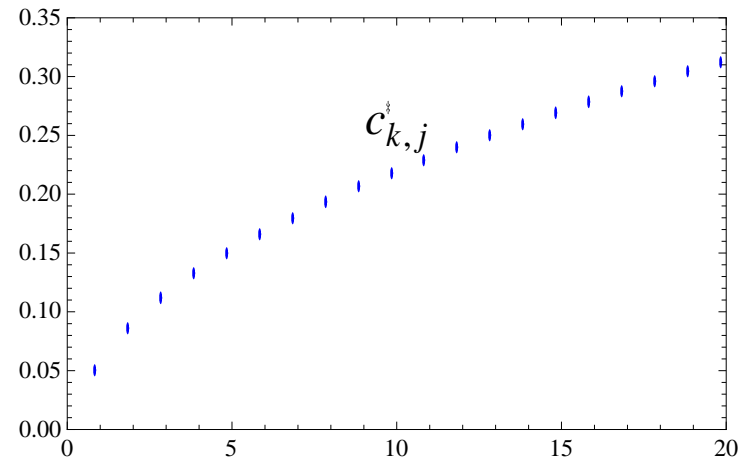$$\mathbf{c}_k = (c_{k1}, \ldots, c_{kk}) \in C_k$$

and serve job $i$ with rate $c_{ki}$

- Assume: Capacity regions $C_k$ compact and symmetric

# Optimal scheduling problem (transient system without arrivals)

- Assume that there are $n$ jobs in the system at time $0$

- What is the optimal way to make the system empty?

- Objective: Minimize the mean delay (or flow time)

- Define: Flow time for policy $\phi$

$$T^{\phi} = \sum_{i=1}^{n} t_i^{\phi}$$

where $t_i$ is the completion time of job $i$

- Define: Operating policies

$$\Phi_n = \{\phi = (\mathbf{c}_1, \ldots, \mathbf{c}_n) : \mathbf{c}_k \in C_k \text{ for all } k\}$$

# Trivial case: One job

- Define:

$$G_1^* = \frac{1}{c_1^*}, \quad c_1^* = \max_{c_1 \in C_1} c_1$$

- Now

$$T^* = \min_{\phi \in \Phi_1} T^\phi = s_1 G_1^*, \quad \phi^* = (\mathbf{c}_1^*)$$

Aalto University
School of Electrical
Engineering

# General case: n jobs

- Define (recursively):

$$G_k^* = \min_{\mathbf{c}_k \in C_k} g_k(\mathbf{c}_k), \quad g_k(\mathbf{c}_k) = \frac{1}{c_{kk}}\left(k - \sum_{i=1}^{k-1} c_{ki} G_i^*\right)$$

- Theorem [Aalto et al. (2011)]: **If**

$$G_1^* < \ldots < G_n^*$$

then

$$T^* = \min_{\phi \in \Phi_n} T^\phi = \sum_{k=1}^n s_k G_k^*, \quad \phi^* = (\mathbf{c}_1^*, \ldots, \mathbf{c}_n^*)$$
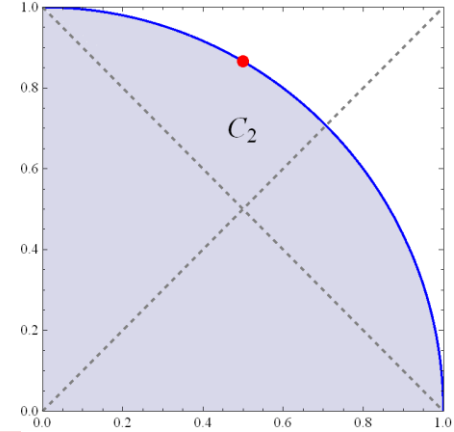
# General case: n jobs (cont.)

- In addition,

$$c^*_{k1} \leq \ldots \leq c^*_{kk} \text{ for all } k$$

- The optimal policy applies the SRPT-FM principle
  - the shortest job is served with the highest rate, etc.
- The optimal rate vector does not depend on the absolute sizes of jobs (only on their order)

Aalto University
School of Electrical
Engineering

11

# Alpha-balls



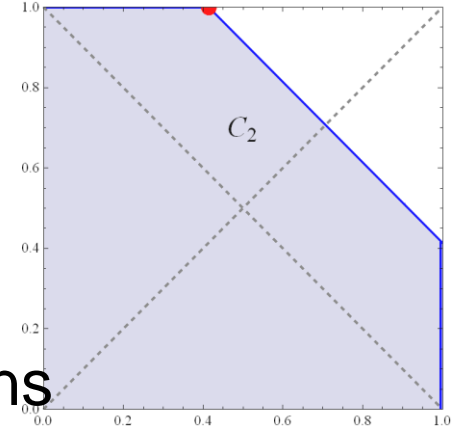- Let $\alpha \geq 1$ and consider capacity regions

$$C_k = \{\mathbf{c}_k \geq 0 : \sum_{j=1}^{k} c_{kj}^{\alpha} \leq 1\}$$

- Now

$$G_k^* = \left( k^{\frac{\alpha}{\alpha-1}} - (k-1)^{\frac{\alpha}{\alpha-1}} \right)^{\frac{\alpha-1}{\alpha}} \quad \text{(increasing in } k)$$

$$c_{kj}^* = \left( \frac{G_j^*}{k} \right)^{\frac{1}{\alpha-1}} \qquad \text{(increasing in } j)$$

# Symmetric polymatroids



- Let $\gamma_1 < \ldots < \gamma_n$ and consider capacity regions

$$C_k = \{\mathbf{c}_k \geq 0 : \sum_{i \in I} c_{ki} \leq \gamma_{|I|}, I \subset \{1, \ldots, n\}\}$$

- Theorem: If $\gamma_1 > \gamma_2 - \gamma_1 > \ldots > \gamma_n - \gamma_{n-1}$, then

$$G_1^* < \ldots < G_n^* \qquad \text{(increasing in } k)$$

$$c_{kj}^* = \gamma_{k-j+1} - \gamma_{k-j} \qquad \text{(increasing in } j)$$

- Optimality result of Sadiq and de Veciana (2010)

# Open questions

- Is it possible to make the implicit condition explicit?

- Optimal scheduling problem for a dynamic system with random arrivals?

- Other objective functions?

**Aalto University**
**School of Electrical**
**Engineering**

# References

- B. Sadiq and G. de Veciana,
  Balancing SRPT prioritization vs opportunistic gain in wireless systems with flow dynamics, in *ITC-22*, 2010

- S. Aalto, A. Penttinen, P. Lassila and P. Osti,
  On the optimal trade-off between SRPT and opportunistic scheduling, in *ACM SIGMETRICS*, 2011

- S. Aalto, A. Penttinen, P. Lassila and P. Osti,
  Optimal size-based opportunistic scheduler for wireless systems, *Queueing Systems* 72, 2012

# The End