

**Optimal Control of Batch Service Queues  
with  
Finite Service Capacity and General Holding Costs**

Samuli Aalto  
EURANDOM  
Eindhoven

## Background

- **Ph.D. Thesis:** “Studies in Queueing Theory”, University of Helsinki, Finland, 1998
  - [1] S. Aalto (1997) Optimal control of batch service queues with Poisson arrivals and finite service capacity, Dep Mathematics, University of Helsinki
  - [2] S. Aalto (1998) Optimal control of batch service queues with compound Poisson arrivals and finite service capacity, *Math Meth Oper Res*
  - [3] S. Aalto (1998) Characterization of the output rate process for a Markovian storage model, *J Appl Prob*
  - [4] S. Aalto (1998) Output of a multiplexer loaded by heterogeneous on-off sources, *Stoch Models*
- Supervisors: Prof. E. Nummelin and Ph.D. T. Lehtonen

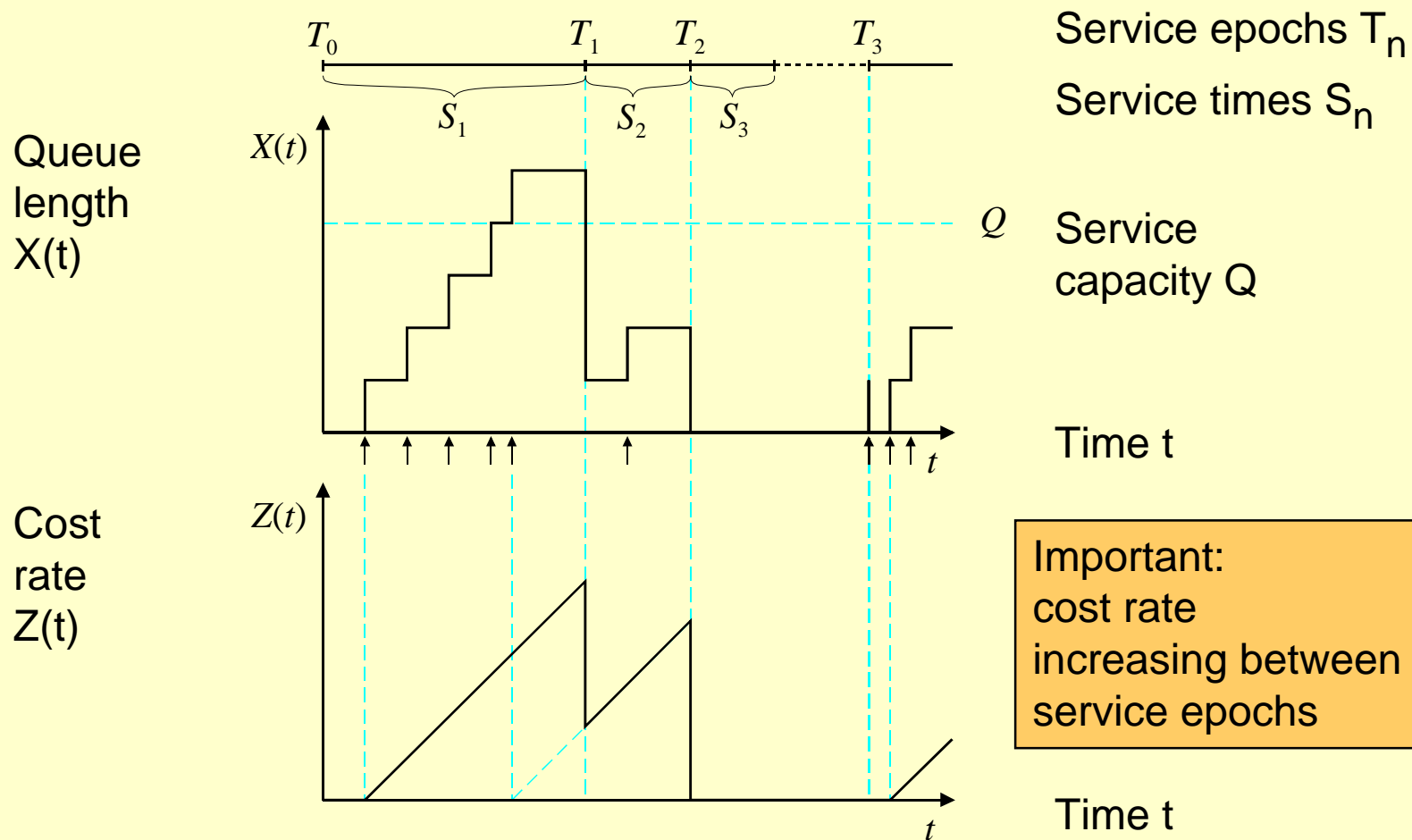
## Contents

- Batch service queue
- Control problem
- Known results
- New results
- Open questions

## Batch service queue

- In an ordinary queue
  - customers are served individually
- In a batch service queue
  - customers are served in batches of varying size
- Additional parameter needed:
  - $Q$  = service capacity = max nr of customers served in a batch

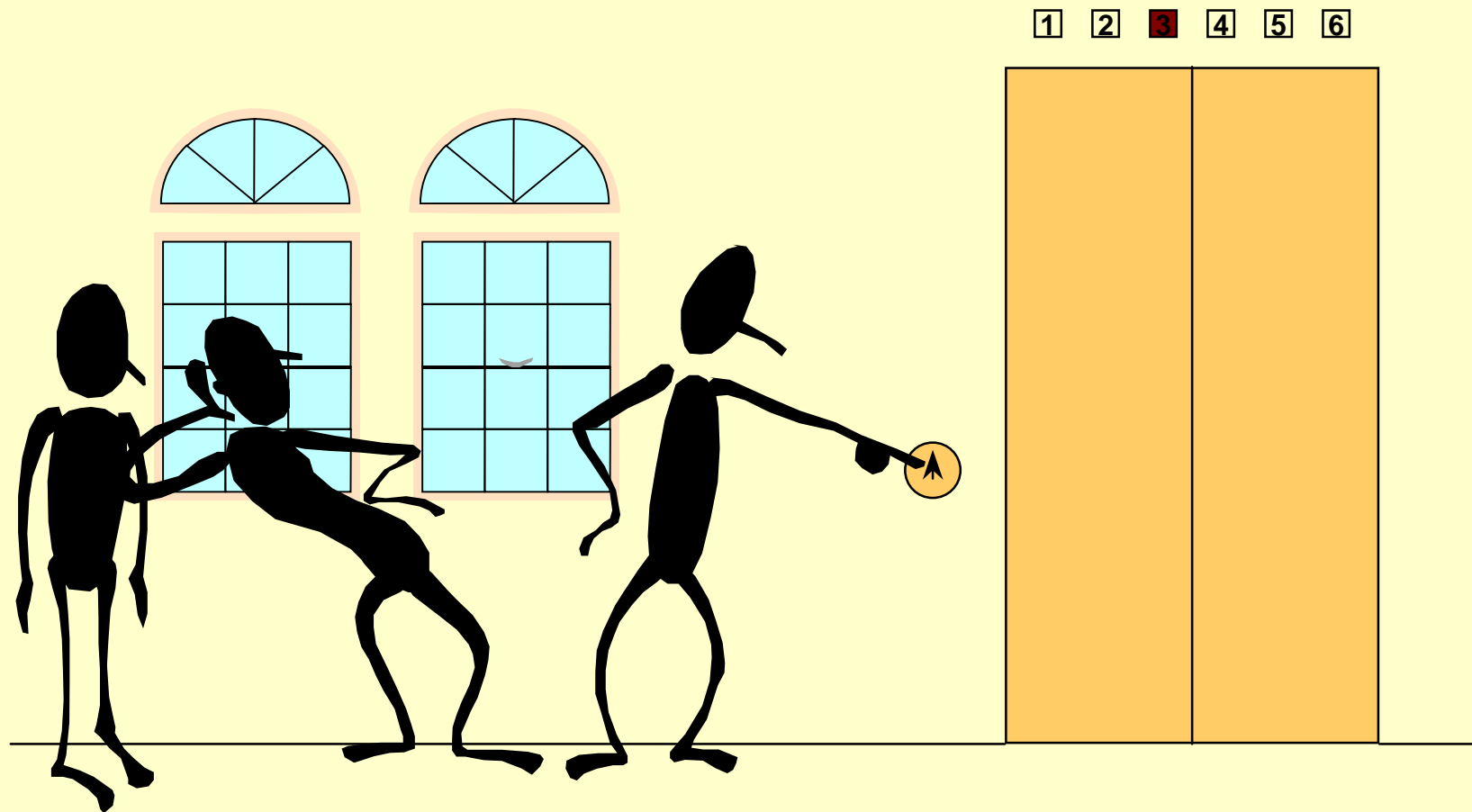
## Evolution



## Queueing models considered

- $M/G(Q)/1$ 
  - Poisson arrivals
  - generally distributed IID service times
  - single server with service capacity  $Q$
- $M^X/G(Q)/1$ 
  - compound Poisson arrivals
  - generally distributed IID service times
  - single server with service capacity  $Q$

## An application



## Contents

- Batch service queue
- Control problem
- Known results
- New results
- Open questions



## Control problem

- Given
  - arrival process  $A(t)$  and
  - service times  $S_n$
- Determine
  - service epochs  $T_n$
  - service batches  $B_n$
- Operating policy  $\pi = ((T_n), (B_n))$ 
  - should be **admissible**

## Optimal control

- Usual operating policy:
  - after a service completion, a new service is initiated as soon as

$$X(t) \geq 1$$

- a service batch includes as many customers as possible
- This is certainly reasonable
- But what is the **optimal** operating policy?
- The answer depends on
  - the **cost structure** and
  - the **objective function**

## Cost structure

- **Holding costs:**  $Z(t)$ 
  - described by the **cost rate** process  $Z(t)$
  - cost rate depends on
    - the nr of waiting customers  $X(t)$ , and
    - the times they have been waiting,  $W_1(t), \dots, W_{X(t)}(t)$
  - called **linear** if

$$Z(t) = h(X(t))$$

- **Serving costs:**  $K + cB_n$ 
  - $K$  per each service batch
  - $c$  per each customer served

## Objective function

- Minimize
  - the **long run average cost**  $\phi^\pi$  or
  - the **discounted cost**  $V_\alpha^\pi$
- Among all the admissible operating policies  $\pi$

## Contents

- Batch service queue
- Control problem
- Known results
- New results
- Open questions

## Known results

	<b>Infinite</b> service capacity $Q = \infty$	<b>Finite</b> service capacity $Q < \infty$
<b>Linear</b> holding costs $z = h(x)$	<b>Case A:</b> - Deb & Serfozo (1973) - Deb (1984)	<b>Case B:</b> - Deb & Serfozo (1973)
<b>General</b> holding costs $z = h(x,w)$	<b>Case C:</b> - Weiss (1979) - Weiss & Pliska (1982)	<b>Case D</b>

## Cases A and B: linear holding costs

A	B

- Deb & Serfozo (1973)
  - Poisson arrivals
  - finite or infinite service capacity
  - average cost & discounted cost
- Deb (1984)
  - compound Poisson arrivals
  - infinite service capacity
  - discounted cost case only
- Result:
  - $h(x)$  is “uniformly increasing”
  - => a **queue length threshold policy** is optimal
- Note: Optimal threshold is never greater than  $Q$

A	B

## Queue length threshold policies

- **Queue length threshold policy**  $\pi_x$  with threshold  $x$ :
  - after a service completion, a new service is initiated as soon as

$$X(t) \geq x$$

- a service batch includes as many customers as possible
- Sufficient to watch over the queue length process  $X(t)$
- Note: the usual operating policy =  $\pi_1$



A	B

## Declaration for cases A and B

- Linear holding costs
  - cost rate remains constant between arrivals
  - system can be reviewed discretely, just when
    - a service has just been completed or
    - the server is free and a new customer arrives
  - Semi-Markov decision technique can be applied
  - no reason to start a new service until the next customer arrives
  - queue length threshold policies are optimal

## Case C: general holding costs & infinite service capacity

C	

- Weiss (1979),  
Weiss & Pliska (1982)
  - compound Poisson arrivals
  - infinite service capacity
  - average cost case only
- Result:
  - $Z(t)$  is increasing (without limits when service is postponed forever)  
=> a **cost rate threshold policy** is optimal

C	

## Cost rate threshold policies

- **Cost rate threshold policy**  $\pi(z)$  with threshold  $z$ :
  - after a service completion, a new service is initiated as soon as

$$Z(t) \geq z$$

- a service batch includes as many customers as possible
  - infinite capacity => all waiting customers
- Sufficient to watch over the cost rate process  $Z(t)$
- If linear and non-decreasing holding costs (  $Z(t) = h(X(t))$  ), then
  - cost rate threshold policies = queue length threshold policies

## Declaration for case C

C	

- Infinite capacity
  - queue can be emptied at every service epoch
  - no reason to watch over the queue length  $X(t)$
  - each service starts a new regeneration cycle (as regards the stationary policies)
- General holding costs
  - system needs to be reviewed continuously
  - Semi-Markov decision technique cannot be applied:
- Cost rate  $Z(t)$  non-decreasing (until the next service)
  - cost rate threshold policies are optimal

## Contents

- Batch service queue
- Control problem
- Known results
- New results
- Open questions

## New results

	<b>Infinite</b> service capacity $Q = \infty$	<b>Finite</b> service capacity $Q < \infty$
<b>Linear</b> holding costs $z = h(x)$	<b>Case A:</b> - Deb & Serfozo (1973) - Deb (1984)	<b>Case B:</b> - Deb & Serfozo (1973)
<b>General</b> holding costs $z = h(x,w)$	<b>Case C:</b> - Weiss (1979) - Weiss & Pliska (1982)	<b>Case D:</b> - Aalto (1997) [1] - Aalto (1998) [2]

## Case D1:

### General holding costs & finite capacity & single arrivals

	D1

- Aalto (1997) [1]
    - Poisson arrivals
    - finite service capacity
    - average cost & discounted cost cases
  - Result:
    - FIFO queueing discipline,
    - **consistent** holding costs and
    - no serving costs included ( $K = c = 0$ )
- => a **cost rate threshold Q-policy** is optimal

	D1

## Consistent holding costs

- Assume that

$$Z(t) = h(X(t), W(t))$$

where  $W(t) = (W_1(t), W_2(t), \dots)$  denotes the vector of

- the waiting times of the customers waiting at time  $t$  (in decreasing order)

- Definition: Holding costs are **consistent** if

$$x \leq x' \text{ and } w \leq w' \quad = \quad h(x, w) \leq h(x', w')$$

- Examples: 1.  $h(x, w) = x$ , 2.  $h(x, w) = w_1 + \dots + w_x$



	D1

## Cost rate threshold Q-policies

- **Cost rate threshold Q-policy**  $\pi_Q(z)$  with threshold  $z$ :
  - after a service completion, a new service is initiated as soon as

$$Z(t) \geq z \quad \text{or} \quad X(t) \geq Q$$

- a service batch includes as many customers as possible
  - finite capacity  $\Rightarrow \min\{X(t), Q\}$
- Necessary to watch over **both** the queue length process  $X(t)$  **and** the cost rate process  $Z(t)$
- If linear and non-decreasing holding costs (  $Z(t) = h(X(t))$  ), then
  - cost rate threshold Q-policies = queue length threshold Q-policies

	D1

## Declaration

- Finite capacity
  - queue **cannot** be emptied at every service epoch
- First key observation:
  - To minimize the holding costs, it is sufficient to consider the class of **Q-policies**
- Second key observation (due to single arrivals):
  - For each Q-policy, the queue becomes empty at every **non-trivial service epoch**
  - such an epoch starts a new regeneration cycle (as regards the stationary Q-policies)

## Q-policies

	D1

- An operating policy  $\pi$  is called a **Q-policy** if
  - after a service completion, a new service is initiated **at latest** when

$$X(t) \geq Q$$

- a service batch includes as many customers as possible
  - finite capacity  $\Rightarrow \min\{X(t), Q\}$

	D1

## Why just Q-policies?

- For each admissible policy  $\pi$ , it is possible to construct such a Q-policy  $\pi^Q$  that

$$X^{\pi^Q}(t) \leq X^{\pi}(t) \quad \forall t$$

- Due to FIFO principle and consistent holding costs, this implies that

$$Z^{\pi^Q}(t) \leq Z^{\pi}(t) \quad \forall t$$

	D1

## Non-trivial service epochs

- Idea: Find such service completions that leave less than  $Q$  customers waiting
- For each Q-policy  $\pi$ , let

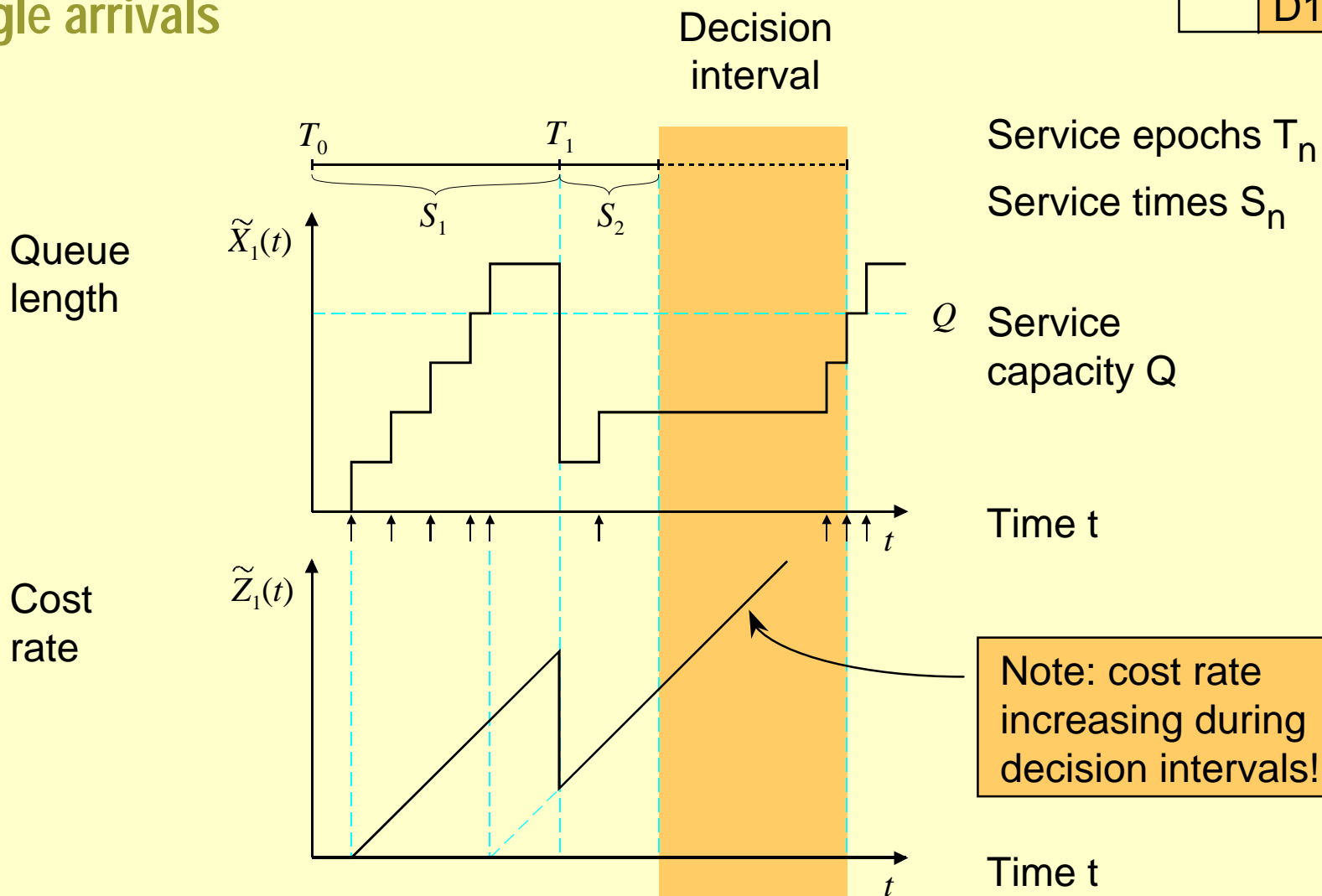
$$N_k^\pi = \min\{n > N_{k-1}^\pi \mid X_n^\pi (T_{n-1}^\pi + S_n) < Q\}$$

- Note:  $N_1$  is the same for all Q-policies  $\pi$
- Definition: **Non-trivial service epochs**

$$\tilde{T}_k^\pi = T_{N_k^\pi}^\pi$$

# Single arrivals

	D1



	D1

## Stationary Q-policies

- Definition: A Q-policy  $\pi$  is called **stationary** if
  - the non-trivial service epochs constitute a renewal sequence, and
  - the process  $(X^\pi, W^\pi)$  is regenerative w.r.t. this sequence
- Long run average cost  $\phi^\pi$ :

$$\phi^\pi = \frac{E[\tilde{C}_1(\tilde{T}_1^\pi)]}{E[\tilde{T}_1^\pi]}$$

where

$$\tilde{C}_1(t) := \int_0^t \tilde{Z}_1(u) du$$

	D1

## Average cost optimal stationary Q-policy

- Denote:

$$\phi = \inf\{\phi^\pi \mid \pi \in \text{stationary Q-policies}\}$$

- Theorem:

- The cost rate threshold Q-policy  $\pi_Q(\phi)$  is average cost optimal among all stationary Q-policies.

- Idea of the proof:

- Show first that the cost rate threshold Q-policy  $\pi_Q(\phi^\pi)$  is better in the average cost sense for any stationary  $\pi$
- Then iterate!



	D1

## Discounted cost optimal stationary Q-policy

- Discounted cost for a stationary Q-policy  $\pi$ :

$$V_{\alpha}^{\pi} = \frac{E[\tilde{D}_{\alpha,1}(\tilde{T}_1^{\pi})]}{1 - E[\exp(-\alpha\tilde{T}_1^{\pi})]}, \quad \text{where } \tilde{D}_{\alpha,1}(t) := \int_0^t e^{-\alpha u} \tilde{Z}_1(u) du$$

- Denote:

$$V_{\alpha} = \inf\{V_{\alpha}^{\pi} \mid \pi \in \text{stationary } Q\text{-policies}\}$$

- Theorem:

– The cost rate threshold Q-policy  $\pi_Q(\alpha V_{\alpha})$  is discounted cost optimal among all stationary Q-policies.

- Similar proof as in the average cost case

	D1

## Discounted cost optimal Q-policy

- Theorem:
  - The cost rate threshold Q-policy  $\pi_Q(\alpha V_\alpha)$  is discounted cost optimal among all Q-policies.
- Idea of the proof:
  - for any policy  $\pi$ , consider a sequence of policies  $\pi_k^*$ , where  $\pi_k^*$  is identical to  $\pi$  up to the  $k^{\text{th}}$  non-trivial service epoch but thereafter changes to the optimal stationary rule
  - the difference between discounted costs of  $\pi$  and  $\pi_k^*$  can be made arbitrarily small by taking  $k$  great enough
  - $\pi_k^*$  is better than  $\pi_{k+1}^*$  in the discounted cost sense
  - $\pi_0^*$  does not depend on the original policy  $\pi$  but is, in fact, the optimal stationary Q-policy  $\pi_Q(\alpha V_\alpha)$

	D1

## Average cost optimal Q-policy

- Theorem:
  - The cost rate threshold Q-policy  $\pi_Q(\phi)$  is average cost optimal among all Q-policies.
- Idea of the proof:
  - Limit of the discounted cost case as  $\alpha$  tends to 0

## Case D2:

### General holding costs & finite capacity & group arrivals

	D2

- Aalto (1998) [2]
    - **compound** Poisson arrivals
    - finite service capacity
    - discounted cost case only
  - Result:
    - FIFO queueing discipline
    - consistent holding costs,
    - no serving costs included ( $K = c = 0$ ) and
    - bounded arrival batches ( $\leq M$ )
- => a **general threshold Q-policy** is optimal

	D2

## General threshold Q-policies

- **General threshold Q-policy**  $\pi_Q(z, \zeta)$  with threshold  $z$  and (non-decreasing) value function  $\zeta$ :
  - after a service completion, a new service is initiated as soon as

$$Z(t) + \zeta(X(t)) \geq z \quad \text{or} \quad X(t) \geq Q$$

- a service batch includes as many customers as possible
  - finite capacity  $\Rightarrow \min\{X(t), Q\}$
- Necessary to watch over **both** the queue length process  $X(t)$  **and** the cost rate process  $Z(t)$
- If linear and non-decreasing holding costs (  $Z(t) = h(X(t))$  ), then
  - general threshold Q-policies = queue length threshold Q-policies

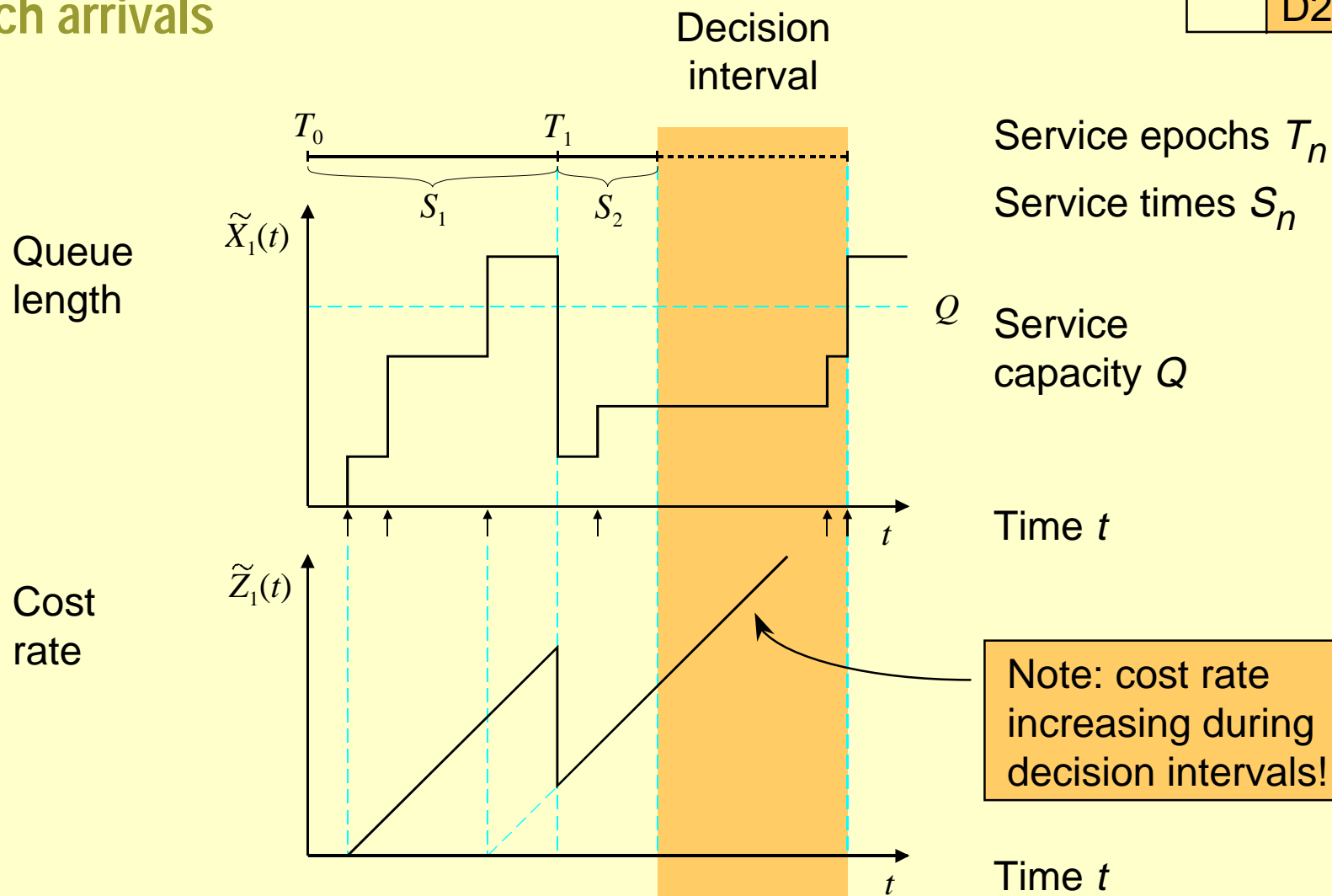
	D2

## Declaration

- Finite capacity
  - queue **cannot** be emptied at every service epoch
- First key observation:
  - To minimize the holding costs, it is (still) sufficient to consider the class of **Q-policies**
- Second key observation (due to FIFO principle):
  - All those customers that remain waiting at a **non-trivial service epoch** arrived at that time => their waiting times are zero
  - such an epoch starts a new “semi-regeneration cycle” (as regards the stationary Q-policies)

## Batch arrivals

	D2



	D2

## Stationary Q-policies

- Definition: A Q-policy  $\pi$  is called **stationary** if
  - the non-trivial service epochs together with the number of customers that remain waiting at those epochs ( $\xi_k^\pi$ ) constitute a Markov renewal sequence, and
  - the process  $(X^\pi, W^\pi)$  is semi-regenerative w.r.t. this sequence
- Discounted cost for a stationary Q-policy  $\pi$ :

$$V^\pi = \mathcal{T}^\pi V^\pi$$

where

$$\mathcal{T}^\pi v(x) := E_x[\tilde{D}_1(\tilde{T}_1^\pi) + \exp(-\tilde{T}_1^\pi)v(\xi_1^\pi)]$$

- Note:  $\xi_k^\pi < M$  for all  $\pi$  and  $k$



	D2

## Discounted cost optimal stationary Q-policy (1)

- Definition: Let

$I_M$  = the non -decreasing functions defined on  $\{0,1,\dots,M - 1\}$

- With each  $v$  in  $I_M$ , associate a general threshold Q-policy  $\pi^v$ :

$$\pi^v = \pi_Q(\alpha v(0), \mathcal{A}^Q_{v^Q})$$

- Proposition: For any  $v$  in  $I_M$  and  $x$  in  $\{0,1,\dots,M-1\}$

$$\mathcal{T}^{\pi^v} v(x) = \inf\{\mathcal{T}^{\pi} v(x) \mid \pi \in Q\text{-policies}\}$$

- Definition: For any  $v$  in  $I_M$  and  $x$  in  $\{0,1,\dots,M-1\}$ , let

$$\tilde{\mathcal{T}}v(x) = \mathcal{T}^{\pi^v} v(x)$$

	D2

## Discounted cost optimal stationary Q-policy (2)

- Definition: Let

$$\mathcal{T}^* = \mathcal{T} \Big|_{I_M^*} \quad \text{where } I_M^* = \{v \in I_M \mid v \leq \mathcal{T}v\}$$

- Proposition:

$$v \in I_M^* \Rightarrow \mathcal{T}v \in I_M^*$$

- Proposition:

$$\mathcal{T}^* \text{ has a unique fixed point } w \in I_M^*$$

- Theorem:

- The general threshold Q-policy  $\pi^w$  is discounted cost optimal among all stationary Q-policies.

	D2

## Discounted cost optimal Q-policy

- Theorem:
  - The general threshold Q-policy  $\pi^w$  is discounted cost optimal among all Q-policies.
- Idea of the proof:
  - for any policy  $\pi$ , consider a sequence of policies  $\pi^*_k$ , where  $\pi^*_k$  is identical to  $\pi$  up to the  $k^{\text{th}}$  non-trivial service epoch but thereafter changes to the optimal stationary rule
  - the difference between discounted costs of  $\pi$  and  $\pi^*_k$  can be made arbitrarily small by taking  $k$  great enough
  - $\pi^*_k$  is better than  $\pi^*_{k+1}$  in the discounted cost sense
  - $\pi^*_0$  does not depend on the original policy  $\pi$  but is, in fact, the optimal stationary Q-policy  $\pi^w$

	B2

## Case B2 (as a special case of D2): Linear holding costs & finite capacity & group arrivals

- Aalto (1998) [2]
  - **compound** Poisson arrivals
  - finite service capacity
  - discounted cost case only
- Corollary (of Case D2):
  - linear holding costs with  $h(x)$  non-decreasing,
  - no serving costs included ( $K = c = 0$ ) and
  - bounded arrival batches

=> a **queue length threshold Q-policy** is optimal

## Contents

- Batch service queue
- Control problem
- Known results
- New results
- Open questions

## Case D1:

### General holding costs & finite capacity & single arrivals

	D1

- If serving costs are included ( $K > 0$ ,  $c > 0$ ),
  - What is the optimal policy in the average cost or discounted cost sense?

## Case D2:

### General holding costs & finite capacity & group arrivals

	D2

- How to get rid of the boundedness assumption concerning the arrival batches?
- If no serving costs are included ( $K = 0, c = 0$ ),
  - Is it true that similar results are valid in the average cost case as in the discounted cost case?
- If serving costs are included ( $K > 0, c > 0$ ),
  - What is the optimal policy in the average cost or discounted cost sense?

## Case B2 (as a special case of D2): Linear holding costs & finite capacity & group arrivals

	B2

- How to get rid of the boundedness assumption concerning the arrival batches?
- If no serving costs are included ( $K = 0, c = 0$ ),
  - Is it true that similar results are valid in the average cost case as in the discounted cost case?
- If serving costs are included ( $K > 0, c > 0$ ),
  - What is the optimal policy in the average cost or discounted cost sense?



