

Task Assignment in a Server Farm with Switching Delays and General Energy-Aware Cost Structure

Esa Hyytiä, Rhonda Righter and Samuli Aalto

Aalto University, Finland
University of California Berkeley, USA

INFORMS APS, July 2013

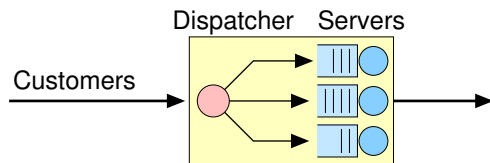


Aalto University
School of Electrical
Engineering

- 1 Task assignment model with switching delay
- 2 Basic dispatching policies
- 3 FPI approach
- 4 Value function for M/G/1-FCFS
 - Switching costs
 - Running costs
 - Holding costs
- 5 Solving the task assignment problem
- 6 Summary of results

Task Assignment with Switching Delay

- **Model**
- **Reference Policies**
- **First Policy Iteration Approach**

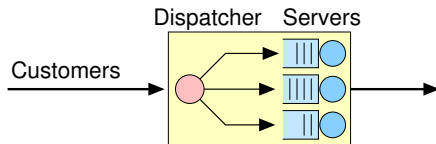


Basic problem $E[T]$

- k parallel queues
- Tasks arrive at rate λ (Poisson process)
- Objective: minimize latency, waiting time, ...

Model for Server Farm

- k parallel servers
- Size-aware setting



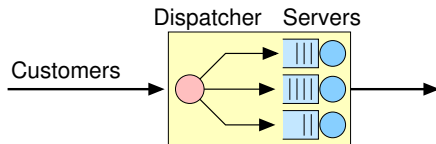
Switching Delays and Energy

Model for Server Farm

- k parallel servers
- Size-aware setting

Distinctive features here

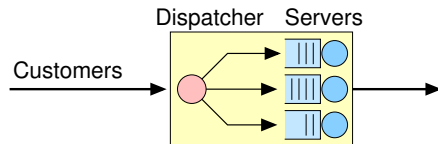
- Idle servers are **switched OFF** to save energy



Switching Delays and Energy

Model for Server Farm

- k parallel servers
- Size-aware setting



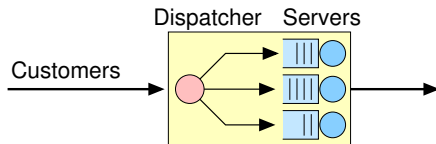
Distinctive features here

- Idle servers are **switched OFF** to save energy
- **Switching ON delay** postpones the start of the service

Switching Delays and Energy

Model for Server Farm

- k parallel servers
- Size-aware setting



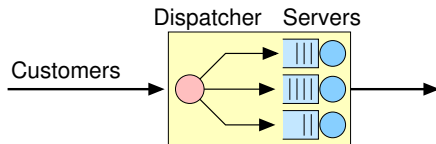
Distinctive features here

- Idle servers are **switched OFF** to save energy
- **Switching ON delay** postpones the start of the service
- **Energy- and Delay-aware** cost structure
 - Switching costs
 - Running costs
 - Holding costs (per job)

Switching Delays and Energy

Model for Server Farm

- k parallel servers
- Size-aware setting



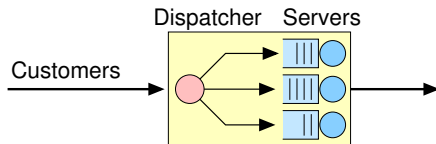
Distinctive features here

- Idle servers are **switched OFF** to save energy
- **Switching ON delay** postpones the start of the service
- **Energy- and Delay-aware** cost structure
 - Switching costs
 - Running costs
 - Holding costs (per job)
- Objective to balance between
 - Energy consumption
 - Performance (e.g., latency)

Switching Delays and Energy

Model for Server Farm

- k parallel servers
- Size-aware setting



Distinctive features here

- Idle servers are **switched OFF** to save energy
- **Switching ON delay** postpones the start of the service
- **Energy- and Delay-aware** cost structure
 - Switching costs
 - Running costs
 - Holding costs (per job)
- Objective to balance between
 - Energy consumption
 - Performance (e.g., latency)

Heterogeneous servers, job-specific costs, . . .

Server Farms with and without Switching Delay

No Switching Delay

Switching Delay

Related models (M/G/1):

- Removable servers, N -policy
 - (Yadin & Naor 1963; Heyman 1968)
 - Service starts when n th customer arrives
- Vacation models, T -policy
 - (Levy & Yechiali 1975; Heyman 1977)
 - Server returns periodically to check the queue
- D -policy, service starts when backlog exceeds d

Related models (M/G/1):

- Removable servers, N -policy
 - (Yadin & Naor 1963; Heyman 1968)
 - Service starts when n th customer arrives
- Vacation models, T -policy
 - (Levy & Yechiali 1975; Heyman 1977)
 - Server returns periodically to check the queue
- D -policy, service starts when backlog exceeds d

Results for switching delay:

- M/G/1 with setup times (Welch, Oper. Res., 1964)
- M/M/ k approximations (Gandhi et al., Sigmetrics'10)
- M/M/ k exact results (Gandhi et al., Sigmetrics'13)

Related models (M/G/1):

- Removable servers, N -policy
 - (Yadin & Naor 1963; Heyman 1968)
 - Service starts when n th customer arrives
- Vacation models, T -policy
 - (Levy & Yechiali 1975; Heyman 1977)
 - Server returns periodically to check the queue
- D -policy, service starts when backlog exceeds d

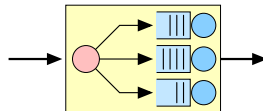
Results for switching delay:

- M/G/1 with setup times (Welch, Oper. Res., 1964)
- M/M/ k approximations (Gandhi et al., Sigmetrics'10)
- M/M/ k exact results (Gandhi et al., Sigmetrics'13)

No delay- and energy-savvy task assignment policies!

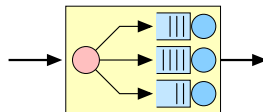
Definition

Static policy chooses the server independently of the queue states



Definition

Static policy chooses the server independently of the queue states

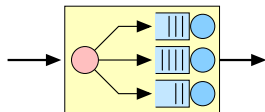


1 Bernoulli splitting (RND)

Choose a queue at random using probabilities p_i

Definition

Static policy chooses the server independently of the queue states



1 Bernoulli splitting (RND)

Choose a queue at random using probabilities p_i

2 Size-Interval-Task-Assignment (SITA)

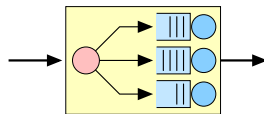
Assignment by the **queue-specific** ranges of job sizes.

“Short jobs to Queue 1 and the rest to Queue 2”

- Proposed in Crovella et. al (Sigmetrics'98) and Harchol-Balter et. al (J. of PDC, 1999)
- Optimal size-aware state-free for FCFS (Feng et. al, -05)

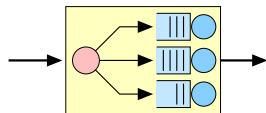
Definition

Actions of a **dynamic policy** depend on the queue states



Definition

Actions of a **dynamic policy** depend on the queue states

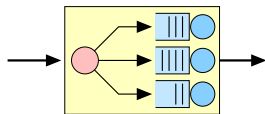


1 Join-the-Shortest-Queue (JSQ)

Optimal when Poisson arrivals, Exp-distributed job sizes, identical servers, and the queue occupancy is known (Haight 1958; Winston 1977)

Definition

Actions of a **dynamic policy** depend on the queue states



1 Join-the-Shortest-Queue (JSQ)

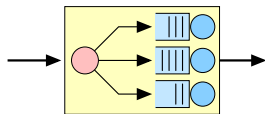
Optimal when Poisson arrivals, Exp-distributed job sizes, identical servers, and the queue occupancy is known (Haight 1958; Winston 1977)

2 Round-robin (RR)

Optimal with identical servers initially in the same state, known routing history and unknown queue occupancy (Ephremides et. al, 1980; Liu&Towsley-94; Liu&Righter -98)

Definition

Actions of a **dynamic policy** depend on the queue states



1 Join-the-Shortest-Queue (JSQ)

Optimal when Poisson arrivals, Exp-distributed job sizes, identical servers, and the queue occupancy is known (Haight 1958; Winston 1977)

2 Round-robin (RR)

Optimal with identical servers initially in the same state, known routing history and unknown queue occupancy (Ephremides et. al, 1980; Liu&Towsley-94; Liu&Righter -98)

3 Least-Work-Left (LWL)

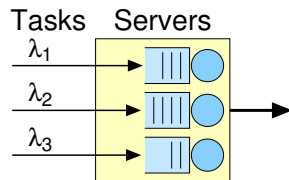
Pick the queue with the shortest backlog (Sharifnia 1997)

- **Decomposition** with a **static basic policy** α
 - Queues receive jobs according to Poisson processes

■ Decomposition with a static basic policy α

- Queues receive jobs according to Poisson processes
- Value function is a sum of the queue-specific value functions

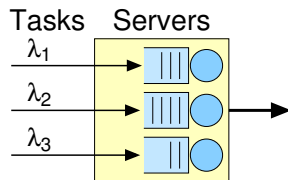
$$v_{\mathbf{z}}(\alpha) = \sum_i v_{\mathbf{z}_i}^{(i)}(\alpha)$$



■ Decomposition with a static basic policy α

- Queues receive jobs according to Poisson processes
- Value function is a sum of the queue-specific value functions

$$v_{\mathbf{z}}(\alpha) = \sum_i v_{\mathbf{z}_i}^{(i)}(\alpha)$$



■ FPI gives a new policy α'

$$\alpha'(x) = \underset{i}{\operatorname{argmin}} \quad h_i(x) + \left(v_{\mathbf{z}_i \oplus x}^{(i)} - v_{\mathbf{z}_i}^{(i)} \right)$$

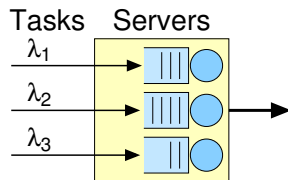
where

- $h_i(x)$ is the *immediate cost* of choosing Queue i for Job x
- $v_{\mathbf{z}_i \oplus x}^{(i)} - v_{\mathbf{z}_i}^{(i)}$ is the mean increase in future costs in Queue i

■ Decomposition with a static basic policy α

- Queues receive jobs according to Poisson processes
- Value function is a sum of the queue-specific value functions

$$v_{\mathbf{z}}(\alpha) = \sum_i v_{z_i}^{(i)}(\alpha)$$



■ FPI gives a new policy α'

$$\alpha'(x) = \underset{i}{\operatorname{argmin}} \quad h_i(x) + \left(v_{\mathbf{z}_i \oplus x}^{(i)} - v_{\mathbf{z}_i}^{(i)} \right)$$

where

- $h_i(x)$ is the *immediate cost* of choosing Queue i for Job x
- $v_{\mathbf{z}_i \oplus x}^{(i)} - v_{\mathbf{z}_i}^{(i)}$ is the mean increase in future costs in Queue i

Idea: The new dynamic policy α' is better than α

Queueing systems

M/M/s

M/M/1

M/G/1-FCFS

M/M/1 & M/M/1/N

M/Cox(r)/1

Krishnan, CDC (1987)

Aalto&Virtamo, NTS-13 (1996)

Sassen et al., Neerlandica (1997)

Koole, CDC (1998)

Bhulai, JAP (2006)

Queueing systems

M/M/s	Krishnan, CDC (1987)
M/M/1	Aalto&Virtamo, NTS-13 (1996)
M/G/1-FCFS	Sassen et al., Neerlandica (1997)
M/M/1 & M/M/1/N	Koole, CDC (1998)
M/Cox(r)/1	Bhulai, JAP (2006)

Size-aware queueing systems

M/G/1 FCFS/LCFS/SRPT	Hyytiä et al., EJOR (2012), Sigmetrics (2012)
M/G/1 (wrt. energy)	Penttinen et al., IPCCC (2011)
M/D/1-PS	Hyytiä et al., ITC (2011)
M/M/1-PS	Hyytiä et al., Performance (2011)

Queueing systems

M/M/s	Krishnan, CDC (1987)
M/M/1	Aalto&Virtamo, NTS-13 (1996)
M/G/1-FCFS	Sassen et al., Neerlandica (1997)
M/M/1 & M/M/1/N	Koole, CDC (1998)
M/Cox(r)/1	Bhulai, JAP (2006)

Size-aware queueing systems

M/G/1 FCFS/LCFS/SRPT	Hyytiä et al., EJOR (2012), Sigmetrics (2012)
M/G/1 (wrt. energy)	Penttinen et al., IPCCC (2011)
M/D/1-PS	Hyytiä et al., ITC (2011)
M/M/1-PS	Hyytiä et al., Performance (2011)

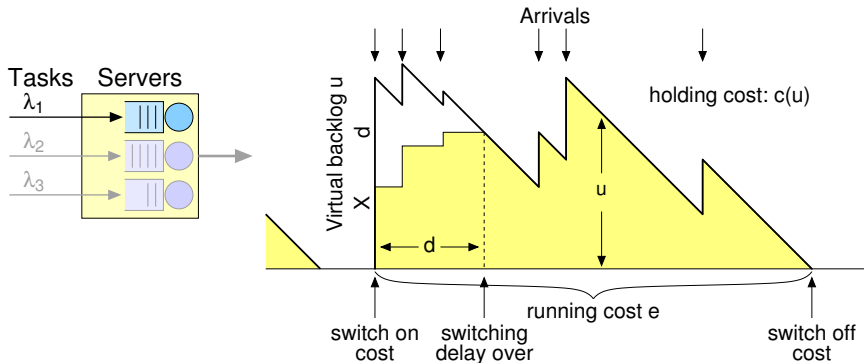
Blocking systems

M/M/s/s	Krishnan, CDC (1986)
M/M/s/k	Leeuwaarden et al. (2001)



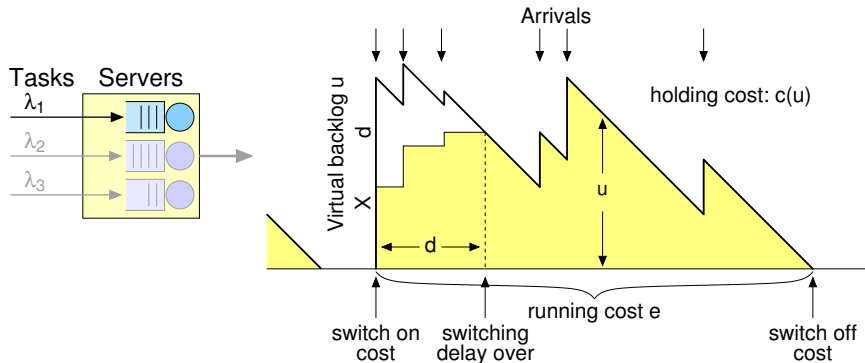
Analysis of Single M/G/1

Cost Structure



¹See (Heyman 1968) and (Feinberg & Kella, 2002) for M/G/1.

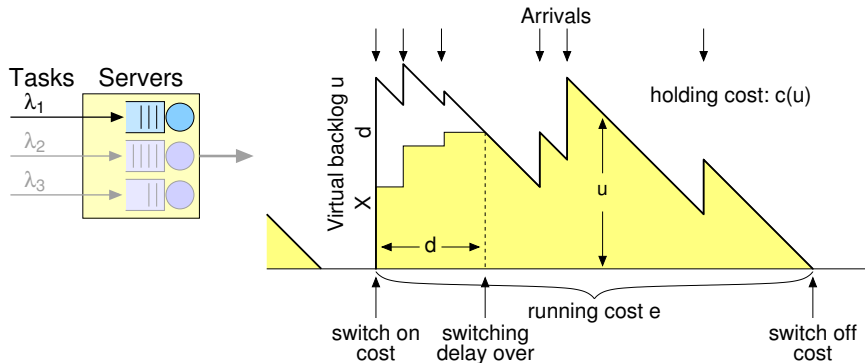
Cost Structure



The queue-specific cost structure¹

(i) **switching costs** ($k_{\text{on}}, k_{\text{off}}$) (per cycle)

¹See (Heyman 1968) and (Feinberg & Kella, 2002) for M/G/1.

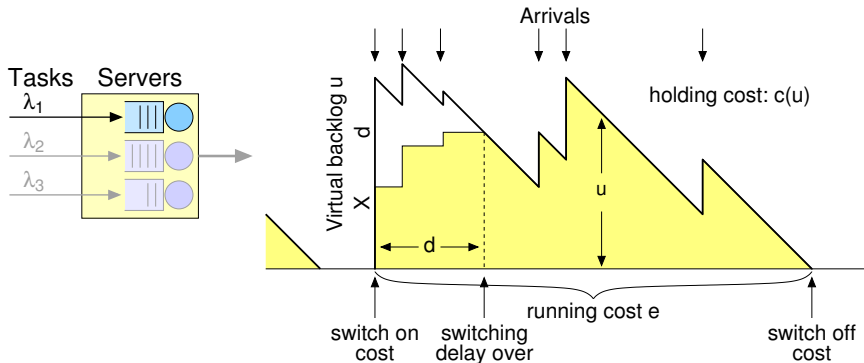


The queue-specific cost structure¹

- (i) **switching costs** ($k_{\text{on}}, k_{\text{off}}$) (per cycle)
- (ii) **running costs** ($e_{\text{on}}, e_{\text{off}}$) (per unit time)

¹See (Heyman 1968) and (Feinberg & Kella, 2002) for M/G/1.

Cost Structure



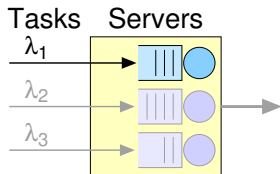
The queue-specific cost structure¹

- (i) **switching costs** ($k_{\text{on}}, k_{\text{off}}$) (per cycle)
- (ii) **running costs** ($e_{\text{on}}, e_{\text{off}}$) (per unit time)
- (iii) **holding cost** $c(u)$ (per unit time), $u = \text{virtual backlog}$

¹See (Heyman 1968) and (Feinberg & Kella, 2002) for M/G/1.



Value Function for M/G/1



■ Formally

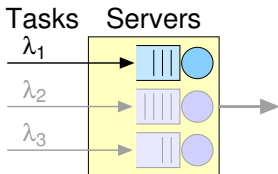
$$v_{\mathbf{z}} \triangleq \lim_{t \rightarrow \infty} E[V_{\mathbf{z}}(t) - r \cdot t]$$

where

- $V_{\mathbf{z}}(t)$ = costs incurred during $(0, t)$ when initially in state \mathbf{z}
- r = long-run mean cost rate



Value Function for M/G/1



Formally

$$v_{\mathbf{z}} \triangleq \lim_{t \rightarrow \infty} E[V_{\mathbf{z}}(t) - r \cdot t]$$

where

- $V_{\mathbf{z}}(t)$ = costs incurred during $(0, t)$ when initially in state \mathbf{z}
- r = long-run mean cost rate
- With FCFS, a sufficient state description is (u, e)
 - u = virtual backlog (measured in time)
 - $e = \begin{cases} 1, & \text{if the server is available} \\ 0, & \text{otherwise (on vacation)} \end{cases}$



Results for M/G/1 without Switching Delay

Cost	Mean rate r_*	Value function $v_*(u) - v_*(0)$
Switching	$\lambda(1 - \rho) \cdot k$	$-\lambda u \cdot k$
Running	$\rho \cdot e$	$u \cdot e$
Holding H_1	$\frac{\lambda E[X^2]}{2(1 - \rho)}$	$\frac{u^2}{2(1 - \rho)}$

Results for M/G/1 with Switching Delay

Cost	Mean rate r_*	Value function $v_*(u) - v_*(0)$
Switching	$\frac{\lambda(1-\rho)}{1+\lambda d} \cdot k$	$-\frac{\lambda u}{1+\lambda d} \cdot k$
Running	$\frac{\rho + \lambda d}{1+\lambda d} \cdot e$	$\frac{u}{1+\lambda d} \cdot e$
Holding H_1	$\frac{\lambda E[X^2]}{2(1-\rho)} + \frac{d(2\rho + \lambda d)}{2(1+\lambda d)}$	$\frac{u^2}{2(1-\rho)} - \frac{d(2\rho + \lambda d) \cdot u}{2(1-\rho)(1+\lambda d)}$

Results for M/G/1 with Switching Delay

Cost	Mean rate r_*	Value function $v_*(u) - v_*(0)$
Switching	$\frac{\lambda(1-\rho)}{1+\lambda d} \cdot k$	$-\frac{\lambda u}{1+\lambda d} \cdot k$
Running	$\frac{\rho + \lambda d}{1 + \lambda d} \cdot e$	$\frac{u}{1 + \lambda d} \cdot e$
Holding H_1	$\frac{\lambda E[X^2]}{2(1-\rho)} + \frac{d(2\rho + \lambda d)}{2(1 + \lambda d)}$	$\frac{u^2}{2(1-\rho)} - \frac{d(2\rho + \lambda d) \cdot u}{2(1-\rho)(1 + \lambda d)}$

Note

- Holding cost with $d = 0$ is the Pollaczek-Khinchine formula
- Switching delay shows up as an extra term in r_{H_1} and $v_{H_1}(u)$
 - Extra cost in $v_{H_1}(u)$ due to switching delay $\propto u$
- Decomposition property (Fuhrmann & Cooper, 1985)



Quadratic Holding Costs

- Linear holding cost corresponds to metrics such as
 - **Latency** (i.e., delay, sojourn time, waiting time)
 - **Slowdown** (ratio of the latency to job size, T/X)
 - ... anything that is directly proportional to T



Quadratic Holding Costs

- Linear holding cost corresponds to metrics such as
 - **Latency** (i.e., delay, sojourn time, waiting time)
 - **Slowdown** (ratio of the latency to job size, T/X)
 - ... anything that is directly proportional to T
- Not everything is linear
 - E.g., longer waiting may cause more customer dissatisfaction \Rightarrow **cost rate increases!**

- Linear holding cost corresponds to metrics such as
 - **Latency** (i.e., delay, sojourn time, waiting time)
 - **Slowdown** (ratio of the latency to job size, T/X)
 - ... anything that is directly proportional to T
- Not everything is linear
 - E.g., longer waiting may cause more customer dissatisfaction \Rightarrow **cost rate increases!**
- What about quadratic costs?

Virtual backlog,	cost rate $\propto U(t)^2$
Latency of Job i ,	cost incurred $\propto (T_i)^2$

- Linear holding cost corresponds to metrics such as
 - Latency (i.e., delay, sojourn time, waiting time)
 - Slowdown (ratio of the latency to job size, T/X)
 - ... anything that is directly proportional to T
- Not everything is linear
 - E.g., longer waiting may cause more customer dissatisfaction \Rightarrow cost rate increases!
- What about quadratic costs?

Virtual backlog,	cost rate $\propto U(t)^2$
Latency of Job i ,	cost incurred $\propto (T_i)^2$

Good news: These can be computed too!



Quadratic Holding Costs

The mean holding cost rate is

$$r_{H2} = E[U^2]$$
$$= \frac{3\lambda^2 E[X^2]^2 + 2\lambda(1-\rho) E[X^3]}{6(1-\rho)^2} + \underbrace{\frac{3\rho + \lambda d}{3(1+\lambda d)} d^2 + \frac{\lambda(2+\lambda d) E[X^2]}{2(1-\rho)(1+\lambda d)} d}_{\text{switching delay}}$$

Quadratic Holding Costs

The mean holding cost rate is

$$\begin{aligned} r_{H2} &= E[U^2] \\ &= \frac{3\lambda^2 E[X^2]^2 + 2\lambda(1-\rho) E[X^3]}{6(1-\rho)^2} + \underbrace{\frac{3\rho + \lambda d}{3(1+\lambda d)} d^2 + \frac{\lambda(2+\lambda d) E[X^2]}{2(1-\rho)(1+\lambda d)} d}_{\text{switching delay}} \end{aligned}$$

The corresponding value function is

$$\begin{aligned} v_{H2}(u) - v_{H2}(0) &= \\ &= \frac{1}{3(1-\rho)} u^3 + \frac{\lambda E[X^2]}{2(1-\rho)^2} u^2 - \underbrace{\left(\frac{3\rho + \lambda d}{3(1-\rho)(1+\lambda d)} d^2 + \frac{\lambda(2+\lambda d) E[X^2]}{2(1-\rho)^2(1+\lambda d)} d \right)}_{\text{switching delay}} u \end{aligned}$$

The mean holding cost rate is

$$r_{H2} = E[U^2]$$

$$= \frac{3\lambda^2 E[X^2]^2 + 2\lambda(1-\rho) E[X^3]}{6(1-\rho)^2} + \underbrace{\frac{3\rho + \lambda d}{3(1+\lambda d)} d^2 + \frac{\lambda(2+\lambda d) E[X^2]}{2(1-\rho)(1+\lambda d)} d}_{\text{switching delay}}$$

The corresponding value function is

$$v_{H2}(u) - v_{H2}(0) =$$

$$\frac{1}{3(1-\rho)} u^3 + \frac{\lambda E[X^2]}{2(1-\rho)^2} u^2 - \underbrace{\left(\frac{3\rho + \lambda d}{3(1-\rho)(1+\lambda d)} d^2 + \frac{\lambda(2+\lambda d) E[X^2]}{2(1-\rho)^2(1+\lambda d)} d \right)}_{\text{switching delay}} u$$

- Mean cost rate (cf. PK) and value function resemble each other
- Switching delay appears as extra terms in both
- In value function, cost of switching delay proportional to $-u$



- For an arbitrary cost function $c(u)$

$$c_1 \triangleq E[c(W_1) + \dots + c(W_{N_u})]$$

$$c_2 \triangleq \lambda E\left[\int_0^{B_u} c(U_t) dt\right]$$

PASTA $\Rightarrow c_1 = c_2$

- For an arbitrary cost function $c(u)$

$$c_1 \triangleq \mathbb{E}[c(W_1) + \dots + c(W_{N_u})]$$

$$c_2 \triangleq \lambda \mathbb{E}\left[\int_0^{B_u} c(U_t) dt\right]$$

$$\text{PASTA} \Rightarrow c_1 = c_2$$

- For **waiting time W** and its square

Linear	$v_W(u) - v_W(0) = \lambda \left(v_{H1}(u) - v_{H1}(0) - \frac{du}{1+\lambda d} \right)$
Quadratic	$v_{W2}(u) - v_{W2}(0) = \lambda \left(v_{H2}(u) - v_{H2}(0) - \frac{d^2 u}{1+\lambda d} \right)$

- For an arbitrary cost function $c(u)$

$$c_1 \triangleq E[c(W_1) + \dots + c(W_{N_u})]$$

$$c_2 \triangleq \lambda E\left[\int_0^{B_u} c(U_t) dt\right]$$

$$\text{PASTA} \Rightarrow c_1 = c_2$$

- For **waiting time** W and its square

Linear	$v_W(u) - v_W(0) = \lambda \left(v_{H1}(u) - v_{H1}(0) - \frac{du}{1+\lambda d} \right)$
Quadratic	$v_{W2}(u) - v_{W2}(0) = \lambda \left(v_{H2}(u) - v_{H2}(0) - \frac{d^2 u}{1+\lambda d} \right)$

- For **latency**, $v_T(u) - v_T(0) = v_W(u) - v_W(0)$
Similarly, an expression for $v_{T2}(u)$ can be obtained

So what do we have?

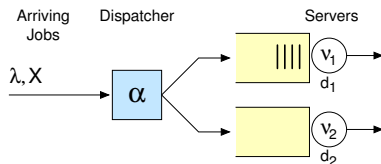
Cost type	mean rate	value function	immediate cost
Switching cost	✓	✓	✓
Running cost	✓	✓	
Waiting time W	✓	✓	✓
Waiting time W^2	✓	✓	✓
Latency T	✓	✓	✓
Latency T^2	✓	✓	✓

⇒ FPI-policy based on a general cost structure!

Solving the Task Assignment Problem

Numerical Examples

Example homogeneous server system

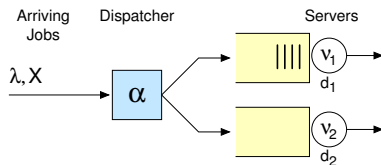


Two homogeneous servers
service rate switching delay

Queue 1: $\nu_1 = 1, d_1 = 1$

Queue 2: $\nu_2 = 1, d_2 = 1$

Example homogeneous server system



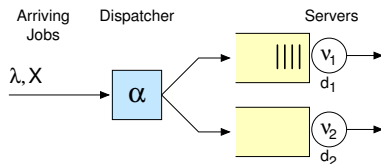
Two homogeneous servers
service rate switching delay

Queue 1: $\nu_1 = 1, \quad d_1 = 1$

Queue 2: $\nu_2 = 1, \quad d_2 = 1$

- $\lambda = 1.5$ and $E[X] = 1$
- Minimize waiting time W

Example homogeneous server system



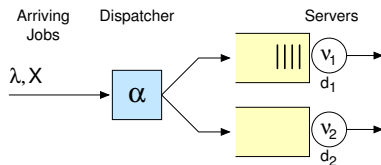
Two homogeneous servers
service rate switching delay

Queue 1: $\nu_1 = 1, \quad d_1 = 1$

Queue 2: $\nu_2 = 1, \quad d_2 = 1$

- $\lambda = 1.5$ and $E[X] = 1$
- Minimize waiting time W
- Basic policy $\alpha = \text{RND}$

Example homogeneous server system



Two homogeneous servers

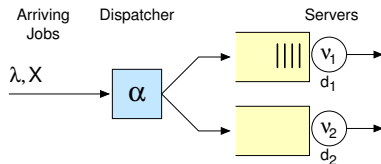
service rate switching delay

Queue 1: $\nu_1 = 1, d_1 = 1$

Queue 2: $\nu_2 = 1, d_2 = 1$

- $\lambda = 1.5$ and $E[X] = 1$
- Minimize waiting time W
- Basic policy $\alpha = \text{RND}$
- Server 1 busy, $u_1 > 0$
Server 2 idle, $u_2 = 0$

Example homogeneous server system

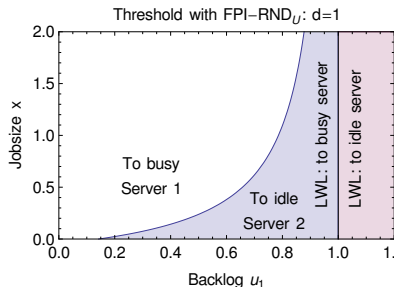


Two homogeneous servers
service rate switching delay

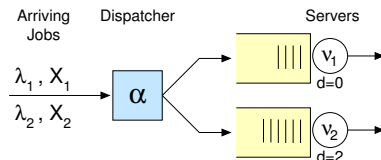
Queue 1: $\nu_1 = 1, d_1 = 1$

Queue 2: $\nu_2 = 1, d_2 = 1$

- $\lambda = 1.5$ and $E[X] = 1$
- Minimize waiting time W
- Basic policy $\alpha = \text{RND}$
- Server 1 busy, $u_1 > 0$
Server 2 idle, $u_2 = 0$
- FPI sends a job to idle Server 2 **earlier than LWL**



Example heterogeneous server system



Two traffic classes:

arrival rate job size

Class 1: $\lambda_1=0.8$ (fixed), $E[X_1]=1$

Class 2: $\lambda_2=0 \dots 0.5$, $E[X_2]=2$

Two heterogeneous servers:

service rate switching delay

Queue 1: $\nu_1 = 1$, $d_1 = 0$ (none)

Queue 2: $\nu_2 = 1$, $d_2 = 2$

Simulation Results

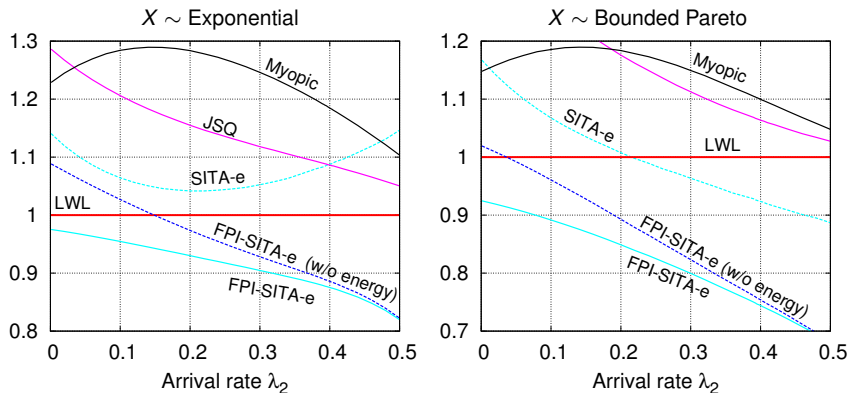


Figure: Mean cost rate with the joint-objective of **waiting time** and **running costs**: $e_1 = 1$ and $e_2 = 2$

- M/G/1 queue with switching delay analyzed

- M/G/1 queue with switching delay analyzed
- Value functions derived with respect to
 - 1 Switching costs [1/cycle]
 - 2 Running costs [1/time]
 - 3 Virtual backlog U_t yielding
 - Waiting time W and its square W^2
 - Latency T and its square T^2

- M/G/1 queue with switching delay analyzed
- Value functions derived with respect to
 - 1 Switching costs [1/cycle]
 - 2 Running costs [1/time]
 - 3 Virtual backlog U_t yielding
 - Waiting time W and its square W^2
 - Latency T and its square T^2
- Enables efficient **task assignment** taking into account
 - Switching delays and service rates
 - Current state of the system
 - Job- and server-specific cost parameters
 - Anticipated future arrivals

- M/G/1 queue with switching delay analyzed
- Value functions derived with respect to
 - 1 Switching costs [1/cycle]
 - 2 Running costs [1/time]
 - 3 Virtual backlog U_t yielding
 - Waiting time W and its square W^2
 - Latency T and its square T^2
- Enables efficient **task assignment** taking into account
 - Switching delays and service rates
 - Current state of the system
 - Job- and server-specific cost parameters
 - Anticipated future arrivals
- Future work: active control of idle servers

- M/G/1 queue with switching delay analyzed
- Value functions derived with respect to
 - 1 Switching costs [1/cycle]
 - 2 Running costs [1/time]
 - 3 Virtual backlog U_t yielding
 - Waiting time W and its square W^2
 - Latency T and its square T^2
- Enables efficient **task assignment** taking into account
 - Switching delays and service rates
 - Current state of the system
 - Job- and server-specific cost parameters
 - Anticipated future arrivals
- Future work: active control of idle servers

Thank you!

- 1 Hyytiä, Righter, Aalto, “*Task Assignment in a Server Farm with Switching Delays and General Energy-Aware Cost Structure*”, submitted, 2013.
- 2 Hyytiä, Aalto, Penttinen, “*Minimizing Slowdown in Heterogeneous Size-Aware Dispatching Systems*”, ACM SIGMETRICS’12.
- 3 Hyytiä, Penttinen, Aalto, “*Size- and State-Aware Dispatching Problem with Queue-Specific Job Sizes*”, EJOR 217(2), 357–370, 2012.

See also,

<http://www.netlab.hut.fi/~esa/dispatching.html>