

State-dependent and Energy-aware Control of Server Farm

Esa Hyytiä, Rhonda Righter and Samuli Aalto

Aalto University, Finland
UC Berkeley, USA

First European Conference on Queueing Theory
ECQT 2014



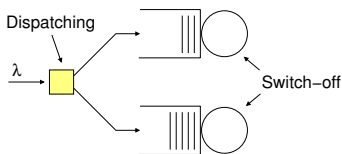
Aalto University
School of Electrical
Engineering

- 1 Model
 - Queueing system with job dispatching
 - Running costs and setup delay
- 2 Optimal static operation
 - Numerical example
- 3 Dynamic operation
 - Value functions for M/G/1
 - Dynamic dispatching and switching off policies
- 4 Conclusions

I. Model

System model:

- n identical FCFS parallel servers
- Jobs dispatched upon arrival
- Running costs at rate e (energy)
- Idle servers can be switched off
- Setup delay of s when switched on
- Objective:



$$\min E[N] + e \cdot E[A]$$

or

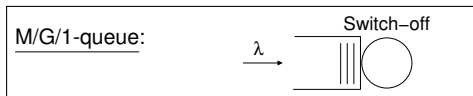
$$\min r_T + r_R$$

where

- $E[N]$ is the mean number in the system ($E[N] = \lambda E[T]$)
- $E[A]$ is the mean number of running servers

Dispatching and Switch-off decisions!

II. Static operation



■ Static switch-off policy:

- 1 **NeverOff**: keep the server always ON
- 2 **InstantOff**: switch off immediately when idle

■ Mean running cost:

$$r_R = \begin{cases} \frac{\lambda(E[X] + s)}{1 + \lambda s} e, & \text{if InstantOff} \\ e, & \text{if NeverOff} \end{cases}$$

■ Mean delay cost:

$$r_T = \begin{cases} \frac{\lambda^2 E[X^2]}{2(1 - \rho)} + \frac{\lambda s(2 + \lambda s)}{2(1 + \lambda s)} + \lambda E[X], & \text{if InstantOff} \\ \frac{\lambda^2 E[X^2]}{2(1 - \rho)} + \lambda E[X], & \text{if NeverOff} \end{cases}$$

The total cost rate under **InstantOff**

$$r_{IO} = \overbrace{\frac{\lambda^2 E[X^2]}{2(1-\rho)} + \frac{\lambda s(2+\lambda s)}{2(1+\lambda s)} + \lambda E[X]}^{\text{Sojourn time}} + \overbrace{\frac{\lambda(E[X] + s)}{1+\lambda s} e}^{\text{Running cost}}$$

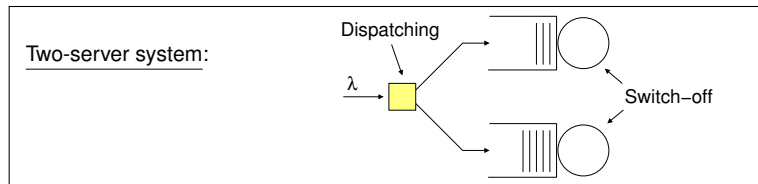
and under **NeverOff**,

$$r_{NO} = \frac{\lambda^2 E[X^2]}{2(1-\rho)} + \lambda E[X] + e$$

Studying $r_{IO} - r_{NO} \Rightarrow$ **InstantOff** better if $e > \frac{\lambda s(2+\lambda s)}{2(1-\rho)}$

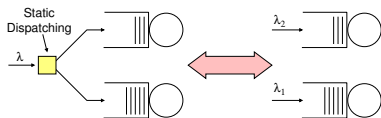
Note: Threshold depends only on $\lambda E[X]$ and λs

Static dispatching: Decomposition



- Static job dispatching
 - *Independent of the queue states*
 - E.g., random split (RND) and SITA¹

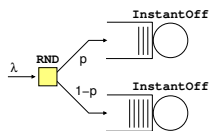
- Decomposition:
Static dispatching:
 $\Rightarrow n$ independent M/G/1 queues



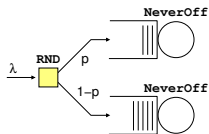
- Mean results available for M/G/1
 - The total cost rate can be computed

¹Size-Interval-Task-Assignment: “short to queue 1, other to queue 2.”

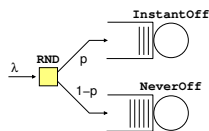
- Two identical servers:
 - Setup time $s = 2$
 - Running cost rate $e = 1$
- Service times $X \sim \text{Exp}(1)$
- Poisson arrival process with rate λ
- RND dispatching (Bernoulli split) w.p. p
- Switch-off policies: **NeverOff** and **InstantOff**



(1) **InstantOff**

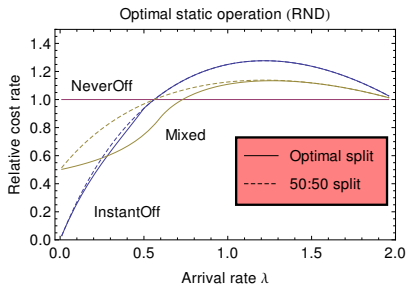


(2) **NeverOff**

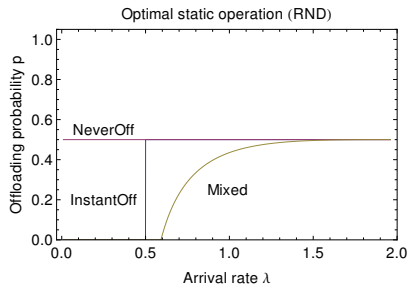


(3) **Mixed**

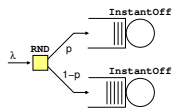
Results



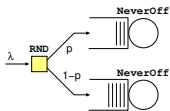
(a) Relative performance



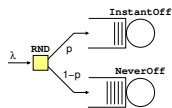
(b) Optimal split



(1) **InstantOff**



(2) **NeverOff**



(3) **Mixed**

Results

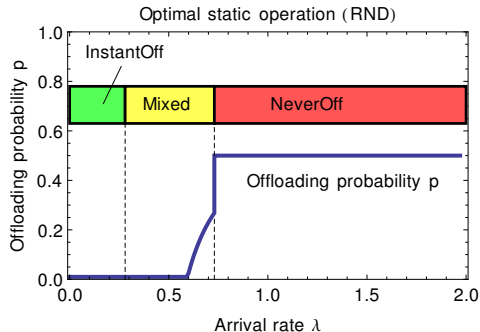


Figure: Optimal operation with RND.

Observations

- Optimal switch-off policy changes as the load increases
InstantOff → **Mixed** → **NeverOff**
- **NeverOff** always splits the jobs uniformly
 - Running costs are fixed, $2 \times e$
 - Uniform split minimizes the mean sojourn time
- **InstantOff** and **Mixed** use
 - Only one server under a very low load
 - Uniform split under a very high load

III. Dynamic operation

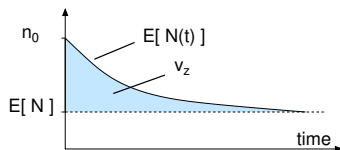
- Dynamic dispatching & switch-off decisions
 - Require state information
 - Can improve the performance
- cf. JSQ vs. RND
- Option to switch-off makes the situation more complicated
- We consider **size- and state-aware** setting

How to capitalize the state information?

Preliminaries

- Consider an arbitrary (stable queueing) system
 - $C_z(t)$ = incurred costs in $(0, t)$ when initially in state z
 - r = mean cost rate
- **Value function** characterizes the expected long-term deviation from the mean cost rate r

$$v_z \triangleq \lim_{t \rightarrow \infty} E[C_z(t) - rt]$$



- For two initial states z_1 and z_2 ,

$$v_{z_2} - v_{z_1}$$

gives the expected difference in the cumulative costs.

Enables the comparison of initial states!

- Virtual backlog u includes the **remaining setup time** δ ,

$$u = \delta + x_1 + \dots, x_n.$$

- Value function w.r.t. **running costs** is²

$$v_R(u) - v_R(0) = \begin{cases} \frac{u}{1+\lambda s} e, & \text{if InstantOff} \\ 0, & \text{if NeverOff} \end{cases}$$

- Value function w.r.t. **sojourn time** is²

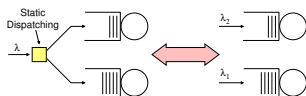
$$v_S(u) - v_S(0) = \begin{cases} \frac{\lambda}{2(1-\rho)} \left(u^2 - \frac{s(2+\lambda s)u}{1+\lambda s} \right) & \text{if InstantOff} \\ \frac{\lambda u^2}{2(1-\rho)} & \text{if NeverOff} \end{cases}$$

The **immediate cost** is equal to the resulting backlog u .

²Hyytiä, Richter, Aalto: Performance Evaluation (2014).

First policy iteration

- Consider a dispatching system with a static policy α_0
- System decomposes to n parallel queues
- Value function is the sum of M/G/1 value functions



$$v_z = \sum_{i=1}^n v^{(i)}(z_i).$$

- Policy iteration step gives a new **dynamic policy**,

$$\alpha(z, x) = \operatorname{argmin}_i \overbrace{c(z_i, x) + v^{(i)}(z_i \oplus x) - v^{(i)}(z_i)}^{\text{Admission cost}}$$

where $c(z_i, x)$ is the immediate cost of server i

1 First policy iteration

(static policy) + (value function) \xRightarrow{FPI} new policy

- Queues are evaluated assuming future jobs according to α_0

2 Lookahead

- Evaluate decisions such as
 - This job to server i
 - Next job to server j (tentatively)
 - Later arriving jobs according to a static α_0
- More accurate evaluation of each possible action
- Yields typically a better policy than FPI

- Two servers

- Server 1: **NeverOff**
Server 2: **InstantOff**

- Setup delay: $s = 2$
- Running cost: $e = 1$

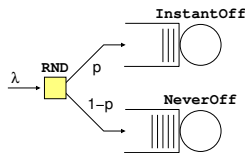
- Objective: Minimize $r_W + r_R$ (waiting time + running costs)

- Reference dispatching policies

- **RND**: random 50:50 split
- **SITA-E**: short jobs to server 1, long to server 2
- **Myopic**: socially optimal if no later arrivals
- **Greedy**: individually optimal choice (only delay)

- Value function based policies

- **FPI**: policy iteration based on SITA-E
- **Lookahead**: “advanced FPI”, considers also the next job (tentatively)



Numerical example: $X \sim \text{Exp}(1)$

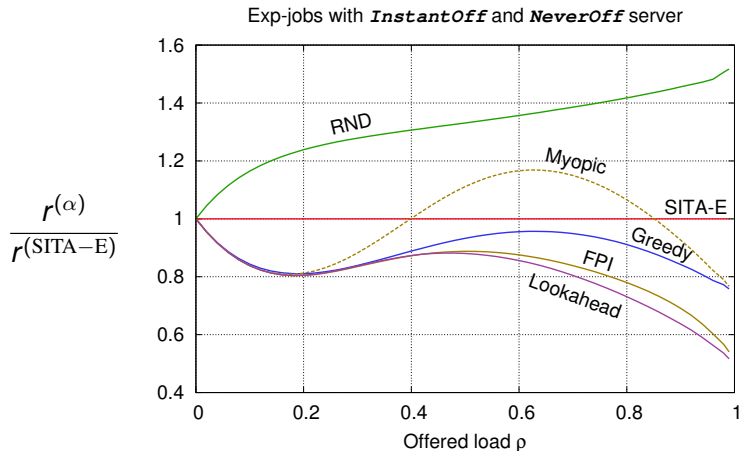


Figure: Relative mean cost rate with the objective of $r_W + r_R$.

Numerical example: $X \sim \text{Pareto}(1)$ (truncated)



Figure: Relative mean cost rate with the objective of $r_W + r_R$.

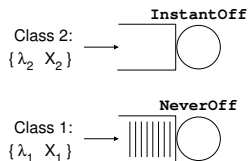


Figure: System according to the static basic policy α_0

Scenario

- Server 1 is `NeverOff` and server 2 `InstantOff`
- Dispatching according to a static α_0
- Server 1 is busy ($u_1 \gg 0$) when server 2 becomes empty

“Should we keep server 2 still running?”

1 Change of roles (renaming)

Idea: swap the roles of the servers?

- Server 1 will receive class 2 jobs, and vice versa
- Server 1 becomes `InstantOff`, and server 2 `NeverOff` (and is thus kept running)

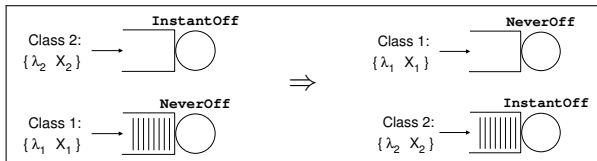


Figure: Swap the roles?

App #2: Improved Switching-off policies (3)

Expected gain can be evaluated with the value functions!

$$\Delta = \frac{u}{2} \left(\left[\frac{\lambda_1}{1 - \rho_1} - \frac{\lambda_2}{1 - \rho_2} \right] u + \frac{\lambda_2 s (2 + \lambda_2 s)}{(1 - \rho_2)(1 + \lambda_2 s)} - \frac{2e}{1 + \lambda_2 s} \right)$$

If $\Delta > 0$, then keep server 2 running

Equivalently,

$$e < \frac{1 + \lambda_2 s}{2} \left(\frac{\lambda_1}{1 - \rho_1} - \frac{\lambda_2}{1 - \rho_2} \right) u + \frac{\lambda_2 s (2 + \lambda_2 s)}{2(1 - \rho_2)}$$

- Uniform RND as α_0 gives

$$e < \frac{\lambda_2 s (2 + \lambda_2 s)}{2(1 - \rho_2)} \quad (\text{same threshold as with a single M/G/1...})$$

- SITA-E as α_0 gives

$$e < \frac{1}{2(1 - \rho_i)} \left((1 + \lambda_2 s)(\lambda_1 - \lambda_2)u + \lambda_2 s (2 + \lambda_2 s) \right)$$

2 Lookahead approach

Similarly as with the dispatching we can ask

A: "Keep server 2 running and assign the next job there?"

For comparison:

B: "Switch off server 2 and assign the next job still there"

C: "Switch off server 2 and assign the next job to server 1"

For each action, we can compute

- The expected costs incurred until the next job arrives
- The expected future costs afterwards (w/ value functions)

While the expected gain is positive, keep server 2 running!

(details omitted)

IV. Conclusions

- Server farm modelled as a queueing system
 - Job dispatching decisions
 - Server switch-off decisions to save energy
 - Setup delay included
- Static control straightforward
 - Mean results available
- Dynamic control is harder
 - Value functions and FPI/Lookahead approaches
 - Can be applied to both dispatching and switching off
- Omitted:
 - Other cost functions (e.g., W^2 , holding costs)
 - Other scheduling disciplines (e.g., LCFS and PS)

Thanks!

- 1 Hyytiä, Righter and Aalto, *Task Assignment in a Heterogeneous Server Farm with Switching Delays and General Energy-Aware Cost Structure*, Performance Evaluation (2014).
- 2 Hyytiä, Righter and Aalto, *Energy-aware Job Assignment in Server Farms with Setup Delays under LCFS and PS*, ITC 2014.