

Load balancing, redundancy, and multi type job and server systems

Urtzi Ayesta

CNRS & Ikerbasque-Univ. Basque Country

Questions/Remarks welcome at : urtzi.ayesta@irit.fr

IFIP Performance 2020, 2/11/2020

Based on:

joint work with:

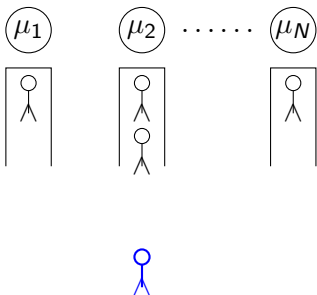
E. Anton, T. Bodas, J.L. Dorsman, M. Jonckheere,
I.M. Verloop

and many other papers by:

Adan, Bonald, Busic, Comte, Gardner,
Harchol-Balter, Hellemans, Hyytiä, Krzesinski,
Mairesse, Moyal, Perry, Righter, Scheller Wolf, van
Houdt, Visschers, Weiss

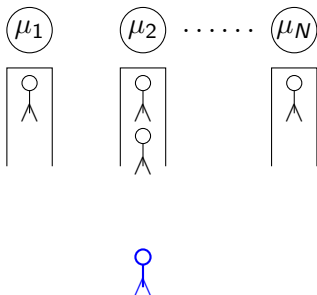
Animations by T. Bodas (Tikz !)

Load balancing



Very active research domain: JSQ, Power of d , Pull-Push based approaches, Jobs with multiple tasks, etc.

Load balancing



Very active research domain: JSQ, Power of d , Pull-Push based approaches, Jobs with multiple tasks, etc.

\implies analysis often approximate or in limiting regimes

Product forms in load balancing systems

- **Markovian queues:** Queues that can be modeled as a *Markov chain*

Product forms in load balancing systems

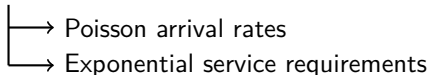
- **Markovian queues:** Queues that can be modeled as a *Markov chain*
 - Poisson arrival rates
 - Exponential service requirements

Product forms in load balancing systems

- **Markovian queues:** Queues that can be modeled as a *Markov chain*
 - Poisson arrival rates
 - Exponential service requirements
- **Goal:** Characterize *stationary distribution* of Markov chain

Product forms in load balancing systems

- **Markovian queues:** Queues that can be modeled as a *Markov chain*

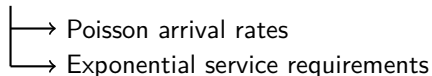


- **Goal:** Characterize *stationary distribution* of Markov chain

Product form

Product forms in load balancing systems

- **Markovian queues:** Queues that can be modeled as a *Markov chain*

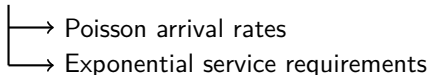


- **Goal:** Characterize *stationary distribution* of Markov chain

└──────────────────┘
↓
Product form \Rightarrow term₁ \times term₂ \times ... term_n

Product forms in load balancing systems

- **Markovian queues:** Queues that can be modeled as a *Markov chain*



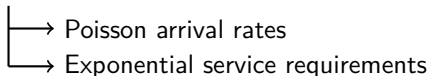
- **Goal:** Characterize *stationary distribution* of Markov chain

Product form

Jackson's 1963 paper on product form queues was considered among Ten Most Influential Titles of Management Sciences First Fifty Years

Product forms in load balancing systems

- **Markovian queues:** Queues that can be modeled as a *Markov chain*



- **Goal:** Characterize *stationary distribution* of Markov chain

Product form

Jackson's 1963 paper on product form queues was considered among Ten Most Influential Titles of Management Sciences First Fifty Years

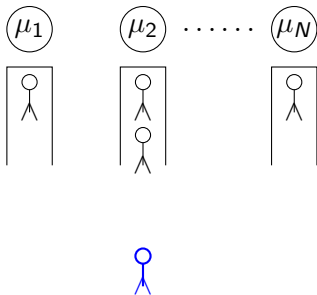
- **Key Idea:** Relate load balancing systems to a central queue architecture

Outline

- ▶ Redundancy
- ▶ Central Queue Architecture
- ▶ Order Independent Descriptor
 - ▶ Redundancy and cancel on complete
- ▶ Aggregated State Descriptor
 - ▶ Redundancy and cancel on start
- ▶ Generalizations:
 - ▶ Token-based framework
 - ▶ Generalized Order Independent
- ▶ Impact of assumptions: scheduling and independence

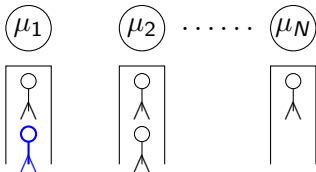
What is redundancy ?

Central idea: create several copies of the same job and use them to minimize latency !



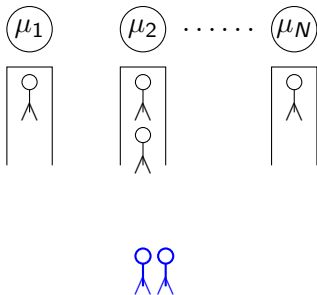
What is redundancy ?

Central idea: create several copies of the same job and use them to minimize latency !



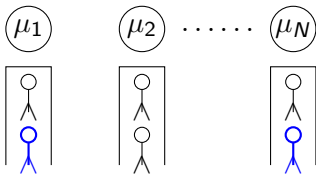
What is redundancy ?

Central idea: create several copies of the same job and use them to minimize latency !



What is redundancy ?

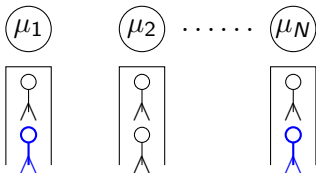
Central idea: create several copies of the same job and use them to minimize latency !



In practice deployed in DNS queries, search engines, Youtube, Mapreduce

What is redundancy ?

Central idea: create several copies of the same job and use them to minimize latency !

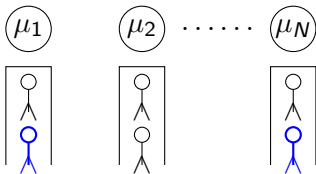


Exploit variability in the workload in different queues !

In practice deployed in DNS queries, search engines, Youtube, Mapreduce

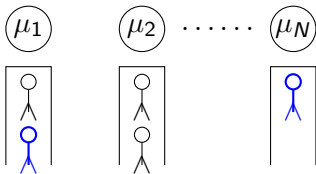
When to cancel the redundant copy ?

A supermarket example



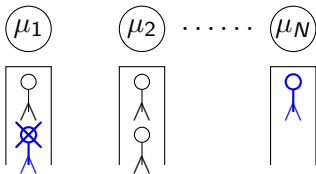
When to cancel the redundant copy ?

A supermarket example



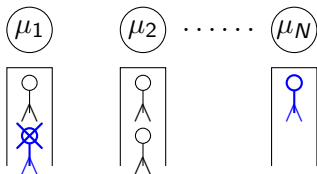
When to cancel the redundant copy ?

A supermarket example



When to cancel the redundant copy ?

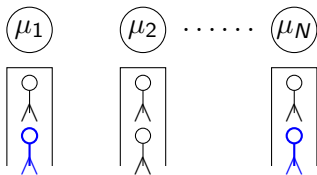
A supermarket example



Cancel on start of service (c.o.s. model)

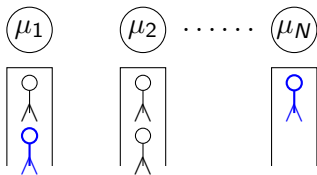
When to cancel the redundant copy ?

The cancel-on-complete variant



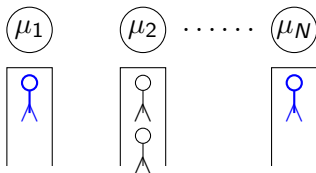
When to cancel the redundant copy ?

The cancel-on-complete variant



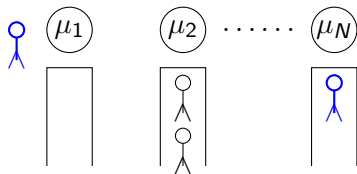
When to cancel the redundant copy ?

The cancel-on-complete variant



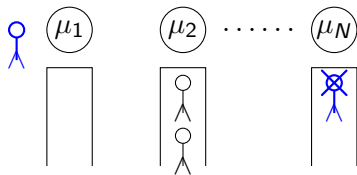
When to cancel the redundant copy ?

The cancel-on-complete variant



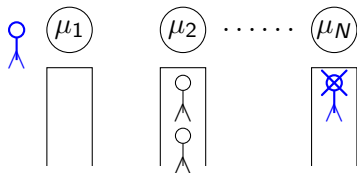
When to cancel the redundant copy ?

The cancel-on-complete variant



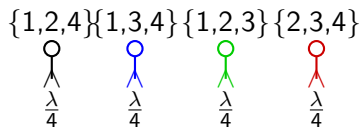
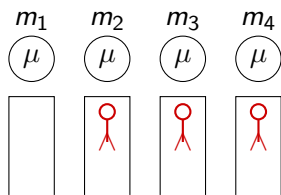
When to cancel the redundant copy ?

The cancel-on-complete variant

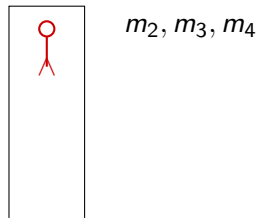
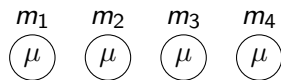
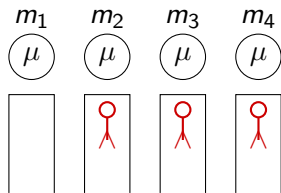


Cancel on completion of service (c.o.c. model)

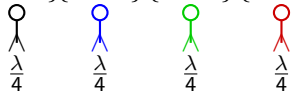
c.o.c. and central queue



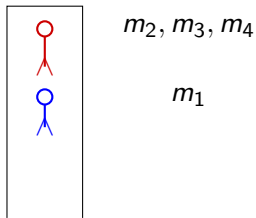
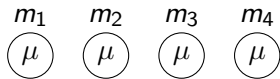
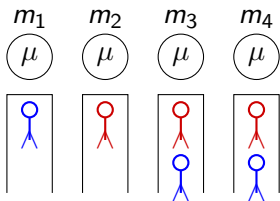
c.o.c. and central queue



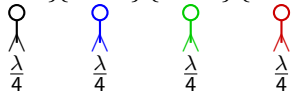
$\{1,2,4\} \{1,3,4\} \{1,2,3\} \{2,3,4\}$



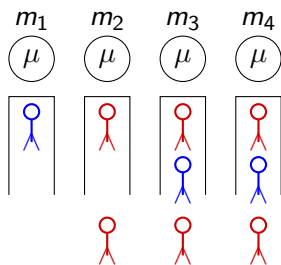
c.o.c. and central queue



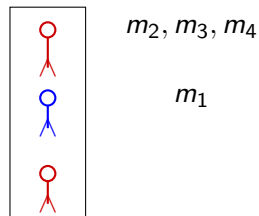
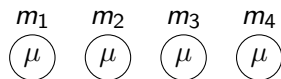
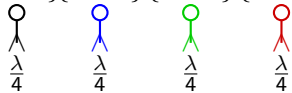
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



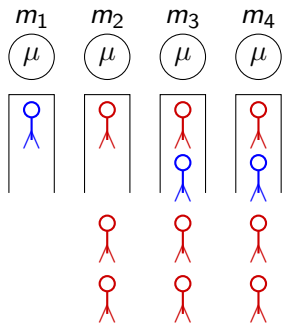
c.o.c. and central queue



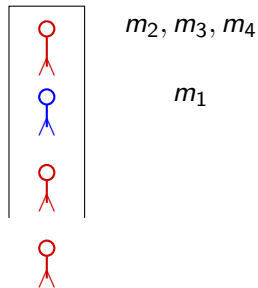
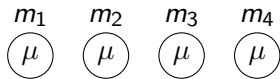
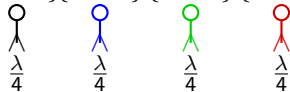
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



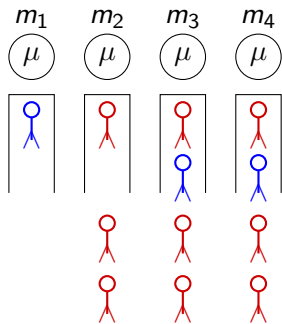
c.o.c. and central queue



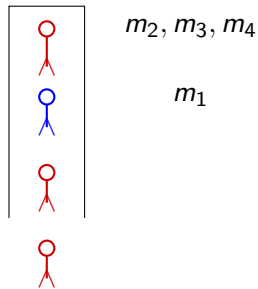
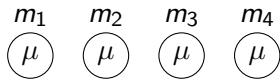
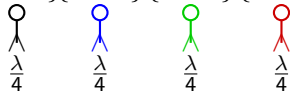
$\{1,2,4\} \{1,3,4\} \{1,2,3\} \{2,3,4\}$



c.o.c. and central queue



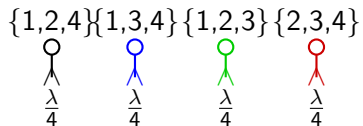
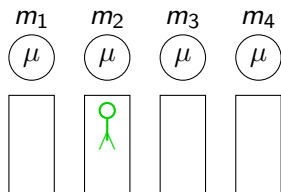
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



Gardner et al. QUESTA 2016

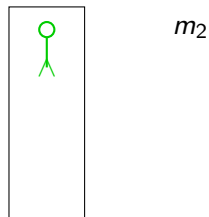
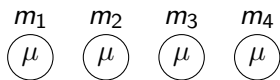
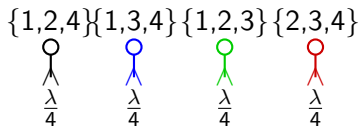
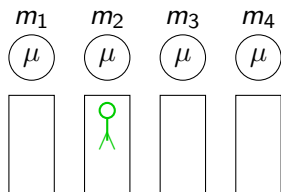
c.o.s. and central queue

$$N = 4 \text{ and } d = 3$$



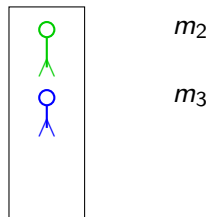
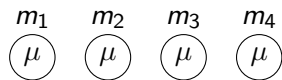
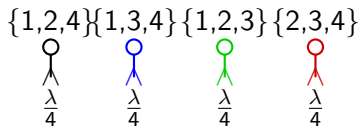
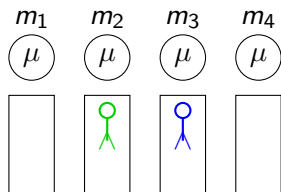
c.o.s. and central queue

$N = 4$ and $d = 3$



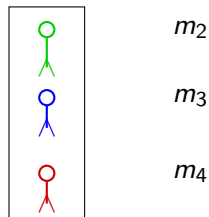
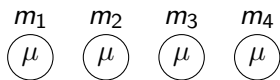
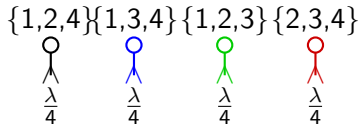
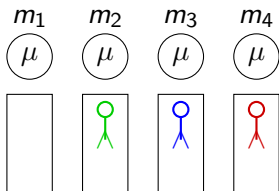
c.o.s. and central queue

$N = 4$ and $d = 3$



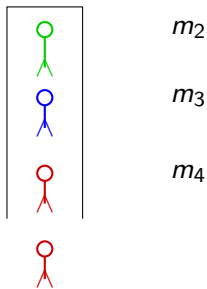
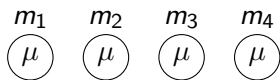
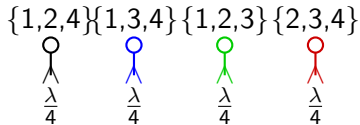
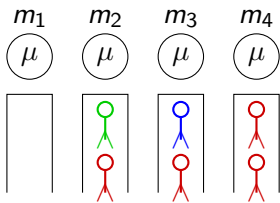
c.o.s. and central queue

$N = 4$ and $d = 3$



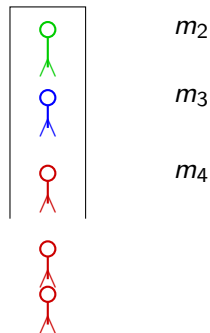
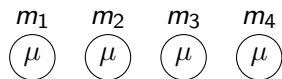
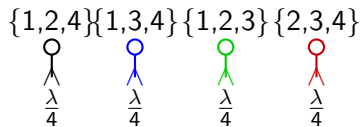
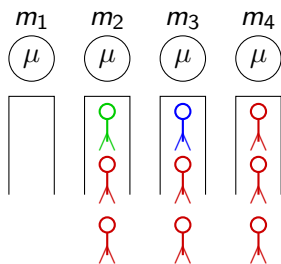
c.o.s. and central queue

$N = 4$ and $d = 3$



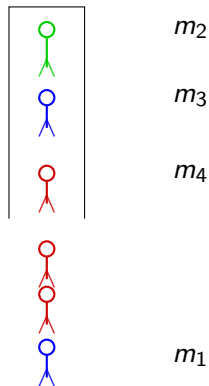
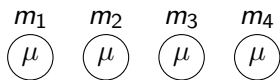
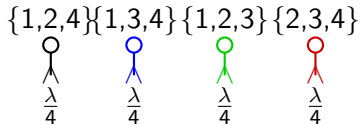
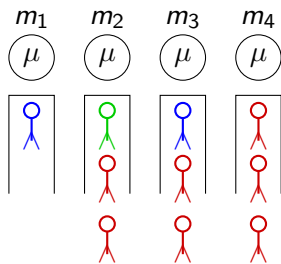
c.o.s. and central queue

$N = 4$ and $d = 3$



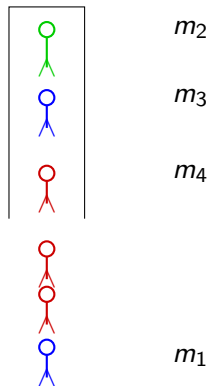
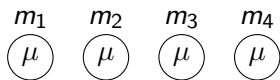
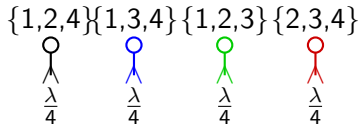
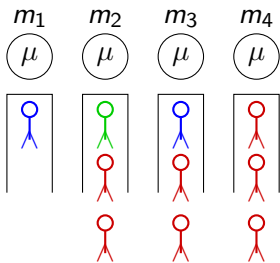
c.o.s. and central queue

$N = 4$ and $d = 3$



c.o.s. and central queue

$N = 4$ and $d = 3$

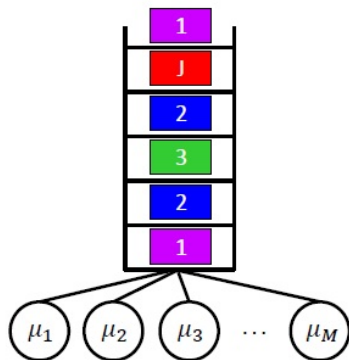


Outline

- ▶ Redundancy
- ▶ Central Queue Architecture
- ▶ Order Independent Descriptor
 - ▶ Redundancy and cancel on complete
- ▶ Aggregated State Descriptor
 - ▶ Redundancy and cancel on start
- ▶ Generalizations:
 - ▶ Token-based framework
 - ▶ Generalized Order Independent
- ▶ Impact of assumptions: scheduling and independence

Central Queue

Central Queue



State descriptors:

- ▶ Job's point of view:
 $c(n) = (c_1, \dots, c_n)$
- ▶ Server's point of view:
 $s = (n_1, M_1, \dots, n_i, M_i)$

Order Independent Queues¹

Classes $i = 1, \dots, N$.

\mathcal{S}_i set of servers that can process class- i

Service is FCFS

State descriptor: $c(n) = (c_1, \dots, c_n)$.

$$\mu(c(n)) = \sum_{s \in \bigcup_{k=1}^n \mathcal{S}_{c_k}} \mu_s$$

¹A. Krzesinski, Order independent queues, in "Queueing Networks: a fundamental approach", Eds: R. Boucherie, N. van Dijk, 2011.

Order Independent Queues¹

Classes $i = 1, \dots, N$.

\mathcal{S}_i set of servers that can process class- i

Service is FCFS

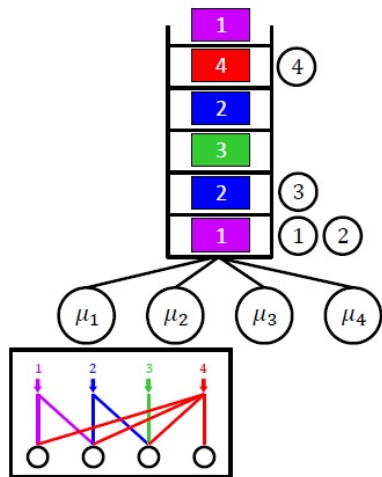
State descriptor: $c(n) = (c_1, \dots, c_n)$.

$$\mu(c(n)) = \sum_{s \in \bigcup_{k=1}^n \mathcal{S}_{c_k}} \mu_s$$

$$\mu(c_1, \dots, c_k) - \mu(c_1, \dots, c_{k-1}) = \sum_{s \in \mathcal{S}_{c_k} \setminus \bigcup_{l=1}^{k-1} \mathcal{S}_{c_l}} \mu_s$$

¹A. Krzesinski, Order independent queues, in "Queueing Networks: a fundamental approach", Eds: R. Boucherie, N. van Dijk, 2011.

Order Independent Queues (example)²

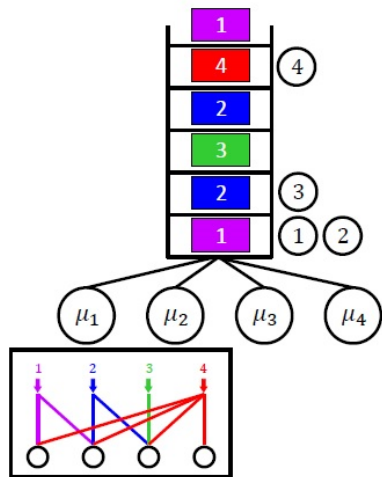


State $c(6) = (1, 2, 3, 2, 4, 1)$

$$\mu(c(6)) = \mu_1 + \mu_2 + \mu_3 + \mu_4$$

²From Gardner and Righter's APS tutorial:
<https://kgardner.people.amherst.edu/>

Order Independent Queues (example)²



State $c(6) = (1, 2, 3, 2, 4, 1)$

$$\mu(c(6)) = \mu_1 + \mu_2 + \mu_3 + \mu_4$$

$$\mu(1, 2) - \mu(1) = \mu_3$$

²From Gardner and Righter's APS tutorial:
<https://kgardner.people.amherst.edu/>

Ol queues

Key Ol properties:

- ▶ For k -th job, its service rate $\mu(c_1, \dots, c_k) - \mu(c_1, \dots, c_{k-1})$ can only depend on what lies ahead.
- ▶ $\mu(c_1, \dots, c_n) = \mu(c_{\sigma(1)}, \dots, c_{\sigma(n)})$

Ol queues

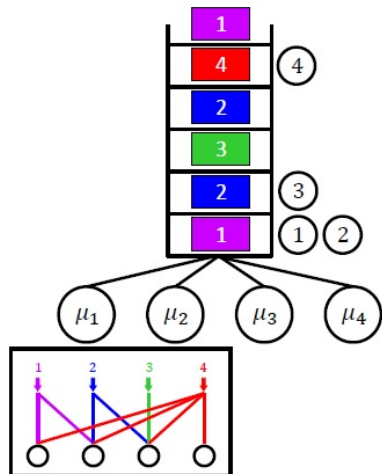
Key Ol properties:

- ▶ For k -th job, its service rate $\mu(c_1, \dots, c_k) - \mu(c_1, \dots, c_{k-1})$ can only depend on what lies ahead.
- ▶ $\mu(c_1, \dots, c_n) = \mu(c_{\sigma(1)}, \dots, c_{\sigma(n)})$

: Theorem: Given Ol properties, the queue is quasi-reversible and the stationary distribution is:

$$\pi(c) = \pi(\emptyset) \prod_{k=1}^n \frac{\lambda_{c_k}}{\mu(c_1, \dots, c_k)}$$

Order Independent Queues (example)



State $c(6) = (1, 2, 3, 2, 4, 1)$

$$\pi(c(6)) = \pi(\emptyset) \left(\frac{\lambda_1}{\mu_{1,2}} \right) \left(\frac{\lambda_2}{\mu_{1,2,3}} \right) \left(\frac{\lambda_3}{\mu_{1,2,3}} \right) \left(\frac{\lambda_2}{\mu_{1,2,3}} \right) \left(\frac{\lambda_4}{\mu} \right) \left(\frac{\lambda_1}{\mu} \right)$$

Sketch Proof of OI

Distribution satisfies partial balance equations:

rate out of $c(n)$ due to departure = rate into $c(n)$ due to
arrival

rate out of $c(n)$ due to class c arrival = rate into $c(n)$
due to class c departure

³T. Bonald, C. Comte, F Mathieu, Performance of Balanced Fairness in Resource Pools: A Recursive Approach, Sigmetrics 2017

Sketch Proof of OI

Distribution satisfies partial balance equations:

rate out of $c(n)$ due to departure = rate into $c(n)$ due to arrival

rate out of $c(n)$ due to class c arrival = rate into $c(n)$ due to class c departure

Quasi reversible \implies Networks of queues have product form stationary distribution

³T. Bonald, C. Comte, F Mathieu, Performance of Balanced Fairness in Resource Pools: A Recursive Approach, Sigmetrics 2017

Sketch Proof of OI

Distribution satisfies partial balance equations:

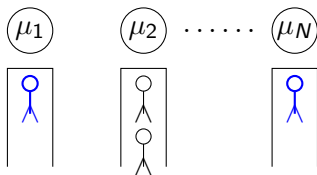
$$\begin{aligned} \text{rate out of } c(n) \text{ due to departure} &= \text{rate into } c(n) \text{ due to arrival} \\ \text{rate out of } c(n) \text{ due to class } c \text{ arrival} &= \text{rate into } c(n) \text{ due to class } c \text{ departure} \end{aligned}$$

Quasi reversible \implies Networks of queues have product form stationary distribution

$\pi(\emptyset)$ can be computed recursively by removing a server s and all classes compatible with that server³

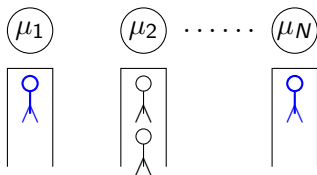
³T. Bonald, C. Comte, F Mathieu, Performance of Balanced Fairness in Resource Pools: A Recursive Approach, Sigmetrics 2017

Redundancy with c.o.c. is OI



$\mu(c_1, \dots, c_n) = \text{sum of rates serving these classes}$

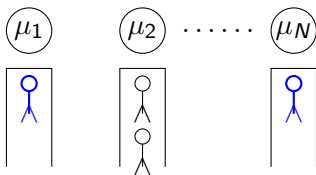
Redundancy with c.o.c. is OI



$\mu(c_1, \dots, c_n)$ = sum of rates serving these classes

$\mu(c_1, \dots, c_n) - \mu(c_1, \dots, c_{n-1})$ = sum of rates of servers available to n class jobs

Redundancy with c.o.c. is OI



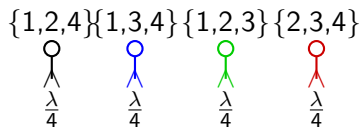
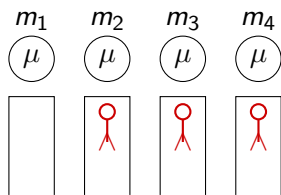
$\mu(c_1, \dots, c_n)$ = sum of rates serving these classes

$\mu(c_1, \dots, c_n) - \mu(c_1, \dots, c_{n-1})$ = sum of rates of servers available to n class jobs

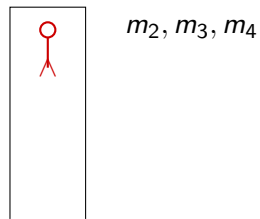
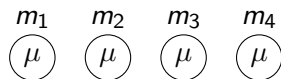
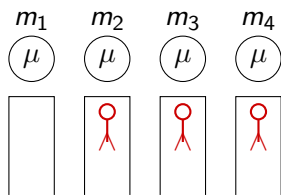
Steady-state distribution is product form!

Key assumption: service time of copies independent!

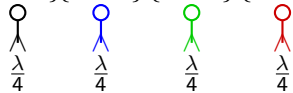
c.o.c. and central queue



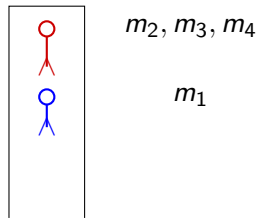
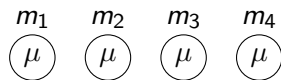
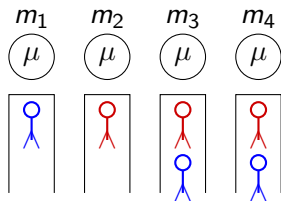
c.o.c. and central queue



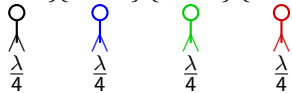
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



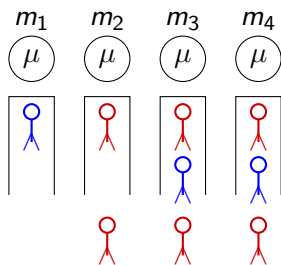
c.o.c. and central queue



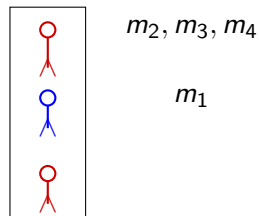
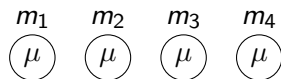
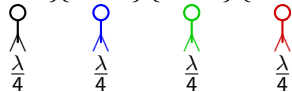
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



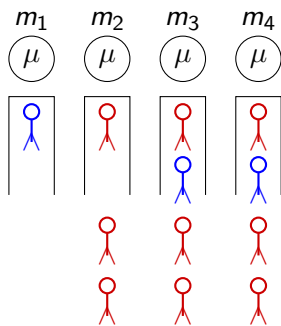
c.o.c. and central queue



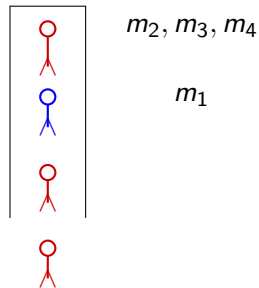
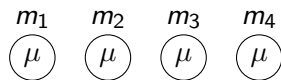
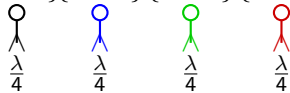
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



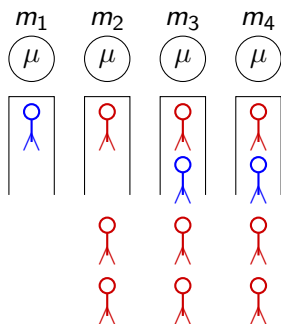
c.o.c. and central queue



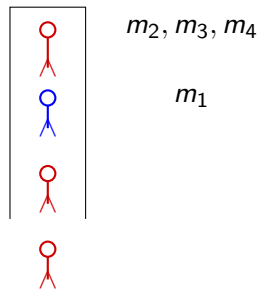
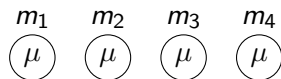
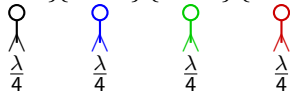
$\{1,2,4\} \{1,3,4\} \{1,2,3\} \{2,3,4\}$



c.o.c. and central queue



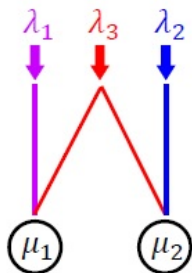
$\{1,2,4\} \{1,3,4\} \{1,2,3\} \{2,3,4\}$



Gardner et al. QUESTA 2016

Performance Evaluation of Redundancy with c.o.c.

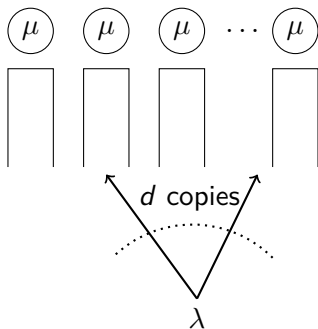
- ▶ Assume particular model: W , N etc.



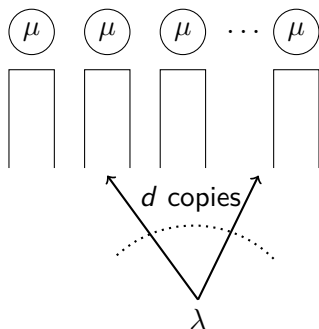
- ▶ Redundancy- d model⁴

⁴Gardner et al, Redundancy- d : The power of d choices for redundancy, OR 2017

redundancy- d with c.o.c.

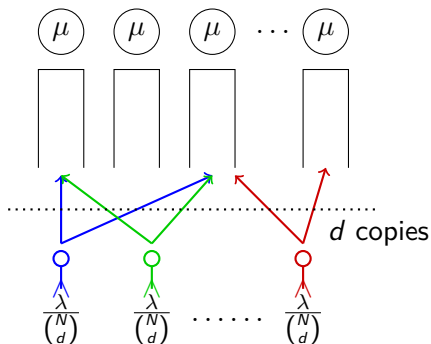


redundancy- d with c.o.c.



- ▶ N homogeneous servers with FIFO discipline
- ▶ Jobs arrivals are Poisson with rate λ
- ▶ Jobs have exponential service requirement
- ▶ Each arrival chooses d servers at random
- ▶ i.i.d. redundant copies for a job placed at d servers

redundancy- d with c.o.c.



- ▶ A special case of the generic multiclass model
- ▶ There are $\binom{N}{d}$ classes
- ▶ Each class has an arrival rate of $\frac{\lambda}{\binom{N}{d}}$

Performance Measures

Mean Delay

$$\mathbb{E}(T^{\text{coc}}) = \sum_{i=d}^k \frac{1}{k\mu \frac{\binom{k-1}{d-1}}{\binom{i-1}{d-1}} - k\lambda}$$

Mean-field limit⁵

$$\mathbb{P}(T^{\text{coc}} > t) = \left(\frac{1}{\rho + (1 - \rho)e^{t\mu(d-1)}} \right)^{\frac{\alpha}{\alpha-1}}$$

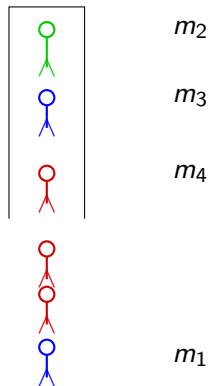
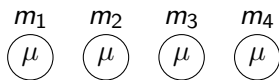
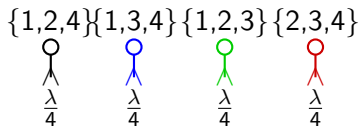
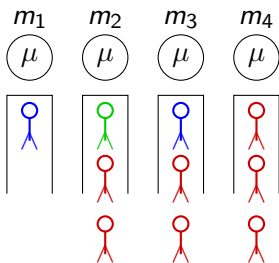
⁵M. Bramson, Yi Lu, B. Prabhakar, Randomized load balancing with general service time distributions. SIGMETRICS 2010: 275-286

Outline

- ▶ Redundancy
- ▶ Central Queue Architecture
- ▶ Order Independent Descriptor
 - ▶ Redundancy and cancel on complete
- ▶ Aggregated State Descriptor
 - ▶ Redundancy and cancel on start
- ▶ Generalizations:
 - ▶ Token-based framework
 - ▶ Generalized Order Independent
- ▶ Impact of assumptions: scheduling and independence

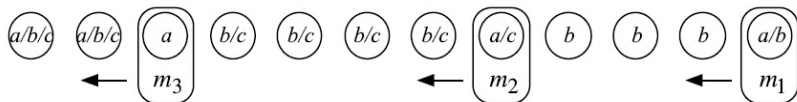
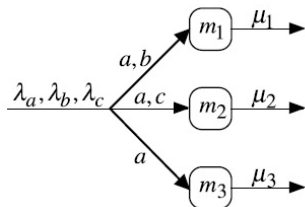
Redundancy c.o.s. is not OI

$N = 4$ and $d = 3$



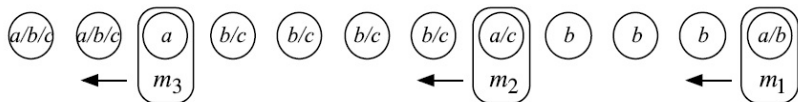
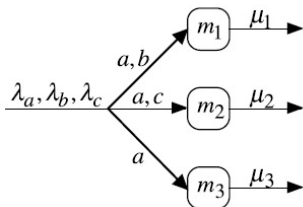
(c_n, \dots, c_1) no product form!

Multi-type job and multi-type servers⁶



⁶Visschers, Adan, Weiss, *A product form solution to a system with multi-type jobs and multi-type servers*, *Queueing Systems*, 2012.

Multi-type job and multi-type servers⁶



- ▶ Markovian descriptor of (aggregated) form
 $s = (n_i, M_i, \dots, n_d, M_d, \dots, M_1)$
- ▶ If i servers are busy, then departure rate is $i\mu$
- ▶ **Assignment rule** determines server in case multiple compatible servers are available

⁶Visschers, Adan, Weiss, *A product form solution to a system with multi-type jobs and multi-type servers*, *Queueing Systems*, 2012.

Multi-type job and multi-type servers (cont.)

Key results: Existence of assignment rule, stability condition, characterization of steady-state distribution

$$\pi(\mathfrak{s}) = \alpha_i^{n_i} \cdots \alpha_1^{n_1} \frac{\Pi_\lambda(\{M_1, \dots, M_i\})}{\Pi_\mu(M_i, \dots, M_1)} \pi(0).$$

Caveat: No efficient way to calculate $\pi(0)$

⁷Adan, Weiss, A skill based parallel service system under FCFS-ALIS: steady state, overloads, and abandonments *Stochastic Systems*, INFORMS, 2014, 4, 250-299

Multi-type job and multi-type servers (cont.)

Key results: Existence of assignment rule, stability condition, characterization of steady-state distribution

$$\pi(\mathfrak{s}) = \alpha_i^{n_i} \cdots \alpha_1^{n_1} \frac{\Pi_\lambda(\{M_1, \dots, M_i\})}{\Pi_\mu(M_i, \dots, M_1)} \pi(0).$$

Caveat: No efficient way to calculate $\pi(0)$

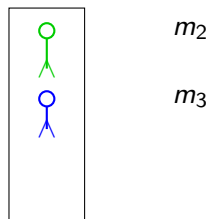
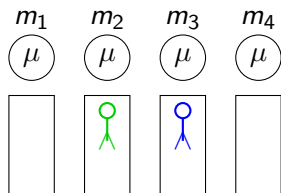
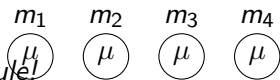
Results of similar nature with [ALIS](#)⁷

⁷Adan, Weiss, A skill based parallel service system under FCFS-ALIS: steady state, overloads, and abandonments *Stochastic Systems*, INFORMS, 2014, 4, 250-299

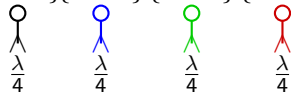
redundancy- d with c.o.s. (example)

$N = 4$ and $d = 3$.

Uniform assignment satisfies assignment rule!



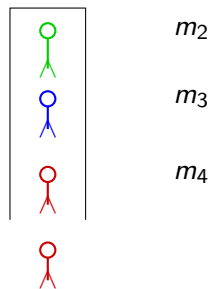
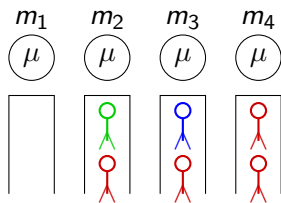
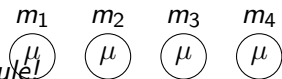
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



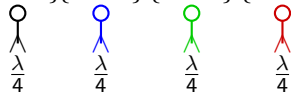
redundancy- d with c.o.s. (example)

$N = 4$ and $d = 3$.

Uniform assignment satisfies assignment rule!



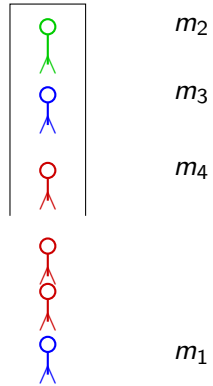
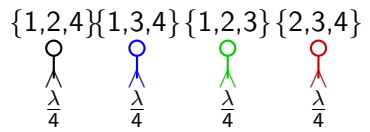
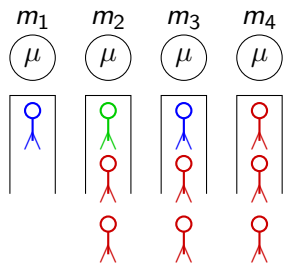
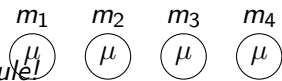
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



redundancy- d with c.o.s. (example)

$N = 4$ and $d = 3$.

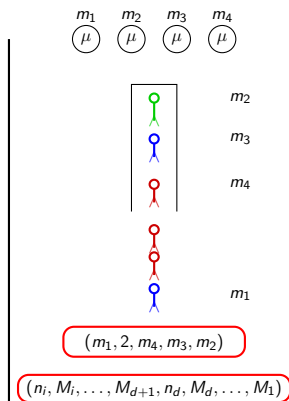
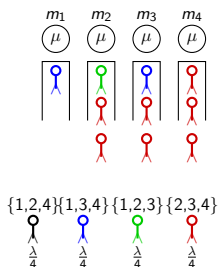
Uniform assignment satisfies assignment rule!



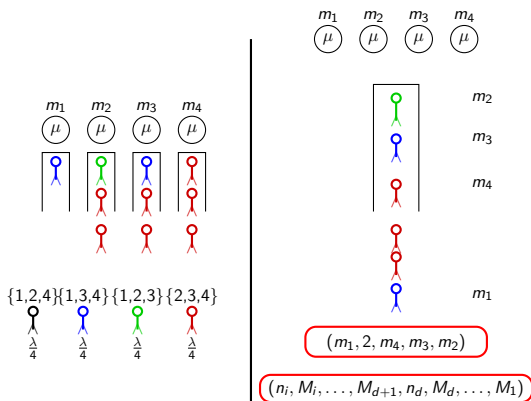
$(m_2, m_3, m_4, 2, m_1)$

$(M_1, \dots, M_d, n_d, M_{d+1}, \dots, M_i, n_i)$

redundancy- d with c.o.s. – results



redundancy- d with c.o.s. – results



For any state $s = (n_i, M_i, \dots, M_{d+1}, n_d, M_d, \dots, M_1)$, we have

$$\pi(s) = r_i^{n_i} \dots r_d^{n_d} \prod_{j=1}^i G_j(K, d) \frac{\pi(0)}{i! \mu^i} \rightarrow \text{product form}$$

redundancy- d with c.o.s. – results⁸

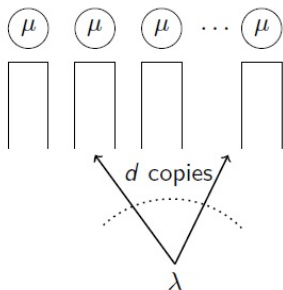
$$\pi(s) = r_i^{n_i} \dots r_d^{n_d} \prod_{j=1}^i G_j(K, d) \frac{\pi(0)}{i! \mu^i}$$

With a bit of algebra, and using the form of $\pi(s)$, we obtain

- ▶ $\pi(0)$, the normalizing constant
- ▶ $p(i)$, probability of i busy servers
- ▶ The P.G.F. of the number of waiting jobs in the system
- ▶ $E(N)$ the expected number of jobs in the system

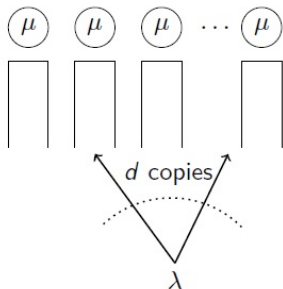
⁸A., Bodas, Verloop, On a unifying product form framework for redundancy models, Performance Evaluation, 2018

Comparing c.o.s. and c.o.c.



If $d = K$, c.o.c. is equivalent to an $M/M/1$ with rate μK
If $d = K$, c.o.s. is equivalent to an $M/M/K$

Comparing c.o.s. and c.o.c.

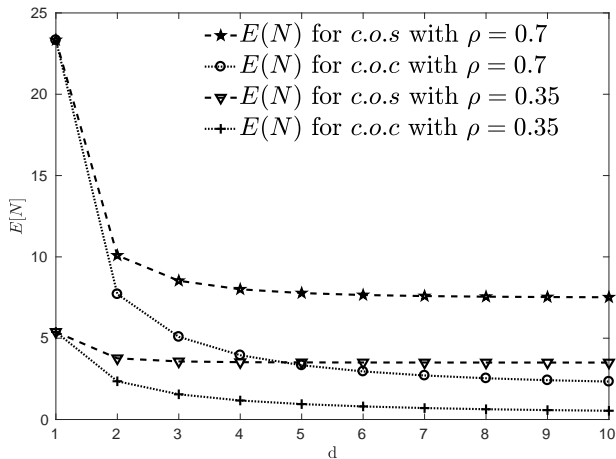


If $d = K$, c.o.c. is equivalent to an $M/M/1$ with rate μK
If $d = K$, c.o.s. is equivalent to an $M/M/K$

► What without i.i.d. assumption?

$d = K, c.o.c. \implies$ single server with rate μ

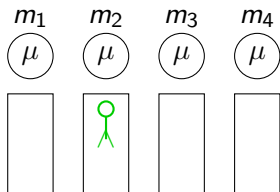
Performance comparison between c.o.s. and c.o.c.



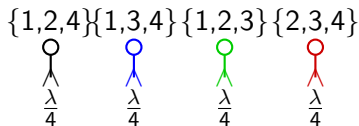
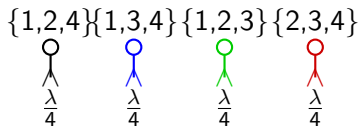
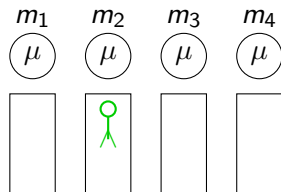
Sample path arguments

$N = 4$ and $d = 3$

redundancy- d with c.o.s.



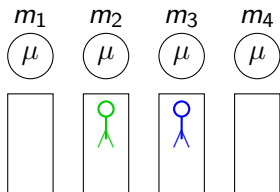
JSW- d



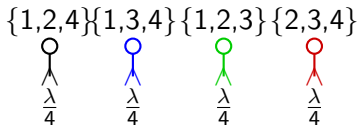
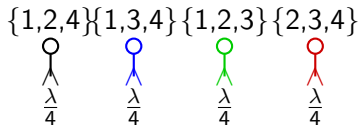
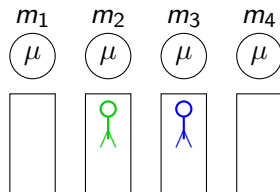
Sample path arguments

$N = 4$ and $d = 3$

redundancy- d with c.o.s.



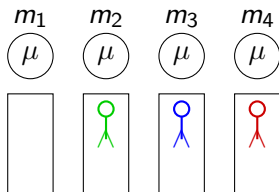
JSW- d



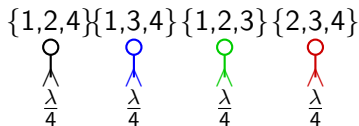
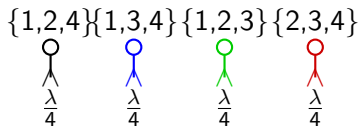
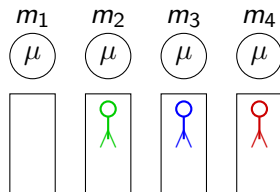
Sample path arguments

$N = 4$ and $d = 3$

redundancy- d with c.o.s.



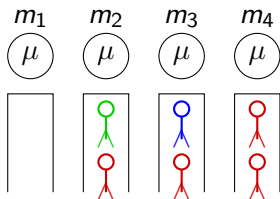
JSW- d



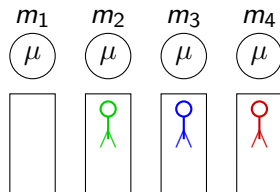
Sample path arguments

$N = 4$ and $d = 3$

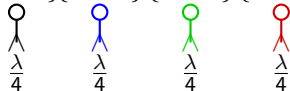
redundancy- d with c.o.s.



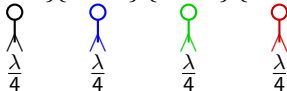
JSW- d



$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



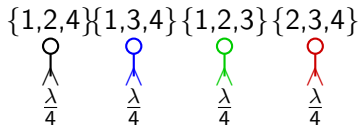
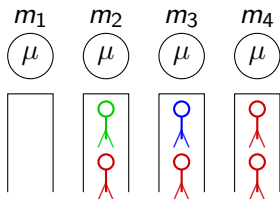
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



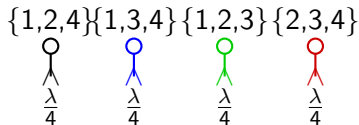
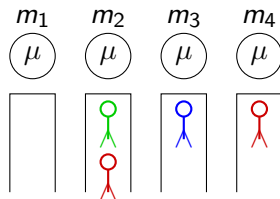
Sample path arguments

$N = 4$ and $d = 3$

redundancy- d with c.o.s.

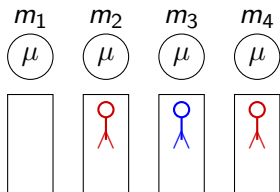


JSW- d

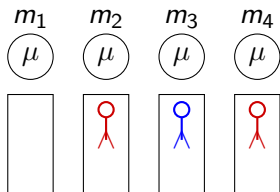


Sample path coupling

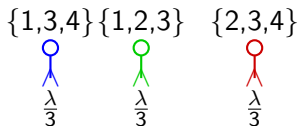
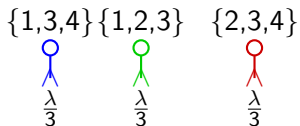
redundancy- d with c.o.s.



JSW- d



Jobs in the two system depart at the same server !



Equivalence with JSW(d)⁹

Proposition

For any given sample-path realization, a given job will be served under both JSW(d) and redundancy-d with c.o.s. in the same server.

We now know $\pi(0)$, $p(i)$, P.G.F for the number of waiting jobs and $E(N)$ for JSW-d through the analysis of Redundancy c.o.s.

⁹A., Bodas, Verloop, On a unifying product form framework for redundancy models, Performance Evaluation, 2018

Redundancy with c.o.s. in mean-field¹⁰

$$\mathbb{P}(T > t) = \left(\lambda^d + (1 - \lambda^d) e^{t(d-1)} \right)^{\frac{1}{d-1}}$$

$$\mathbb{E}(T^{\text{cos}}) = \sum_{n=0}^{\infty} \frac{\lambda^{dn}}{1 + n(d-1)}$$

¹⁰T. Hellemans, B. van Houdt, On the power-of-d-choices with least loaded server selection, Sigmetrics 2018

Redundancy with c.o.s. in mean-field¹⁰

$$\mathbb{P}(T > t) = \left(\lambda^d + (1 - \lambda^d) e^{t(d-1)} \right)^{\frac{1}{d-1}}$$

$$\mathbb{E}(T^{\text{cos}}) = \sum_{n=0}^{\infty} \frac{\lambda^{dn}}{1 + n(d-1)}$$

In the mean-field limit:

$$\mathbb{E}(T^{\text{cos}}) \geq \mathbb{E}(T^{\text{coc}})$$

¹⁰T. Hellemans, B. van Houdt, On the power-of-d-choices with least loaded server selection, Sigmetrics 2018

Outline

- ▶ Redundancy
- ▶ Central Queue Architecture
- ▶ Order Independent Descriptor
 - ▶ Redundancy and cancel on complete
- ▶ Aggregated State Descriptor
 - ▶ Redundancy and cancel on start
- ▶ Generalizations:
 - ▶ Token-based framework
 - ▶ Generalized Order Independent
- ▶ Impact of assumptions: scheduling and independence

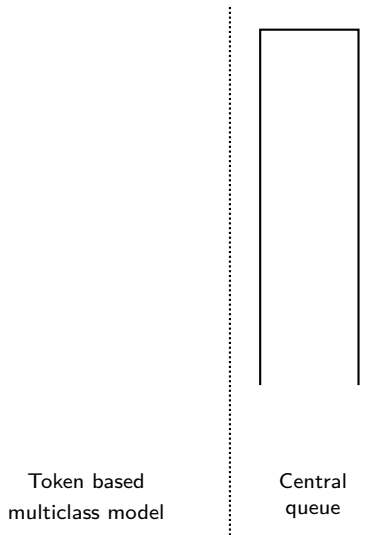
Token based framework¹¹

- ▶ Multiclass jobs and tokens for service
- ▶ Product form stationary distribution
- ▶ Subsumes OI and Visschers et al., and more

¹¹A., Bodas, Dorsman, Verloop, A token-based central queue with order-independent service rates , to appear in OR

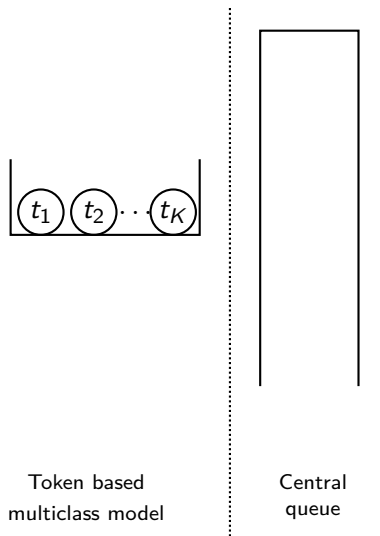
A token based central queue

A token based central queue



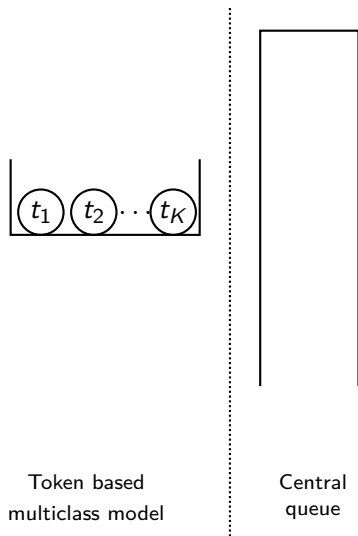
- A single central queue

A token based central queue



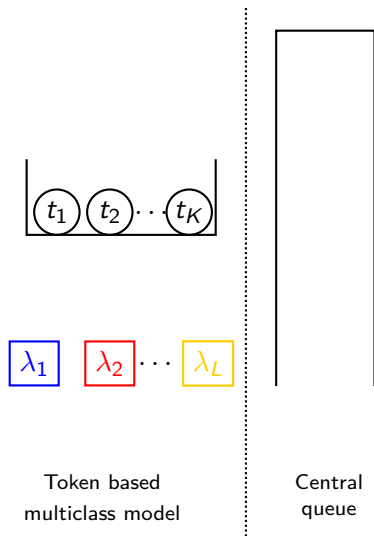
- A single central queue
- A set $\mathcal{T} = \{t_1, \dots, t_K\}$ of K tokens

A token based central queue



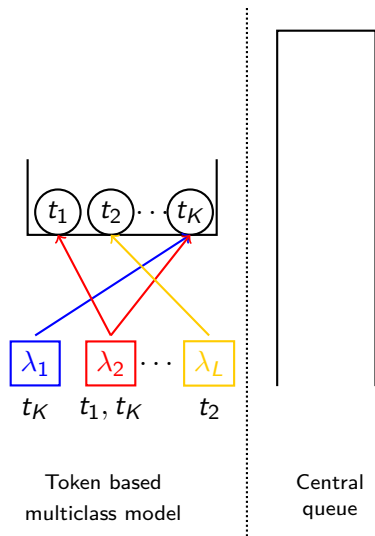
- A single central queue
- A set $\mathcal{T} = \{t_1, \dots, t_K\}$ of K tokens
- Only jobs with tokens are served

A token based central queue



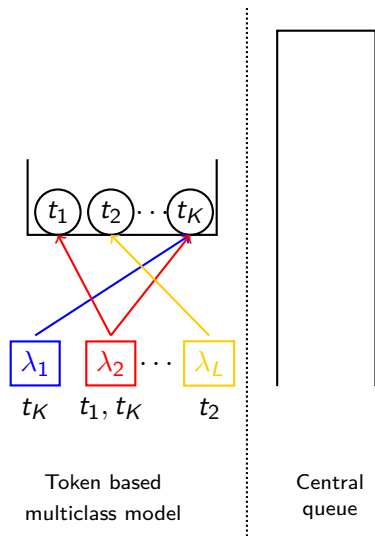
- A single central queue
- A set $\mathcal{T} = \{t_1, \dots, t_K\}$ of K tokens
- Only jobs with tokens are served
- There are L job classes
- Poisson arrival rate λ_i for class i

A token based central queue

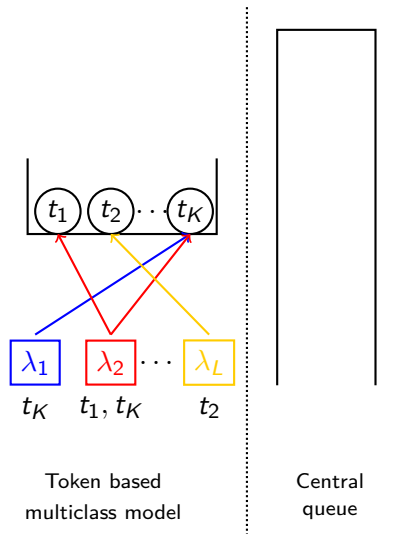


- A single central queue
- A set $\mathcal{T} = \{t_1, \dots, t_K\}$ of K tokens
- Only jobs with tokens are served
- There are L job classes
- Poisson arrival rate λ_i for class i
- A compatibility graph (arbitrary)

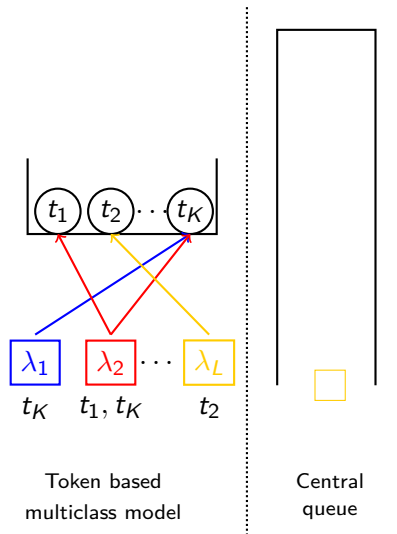
A token based central queue



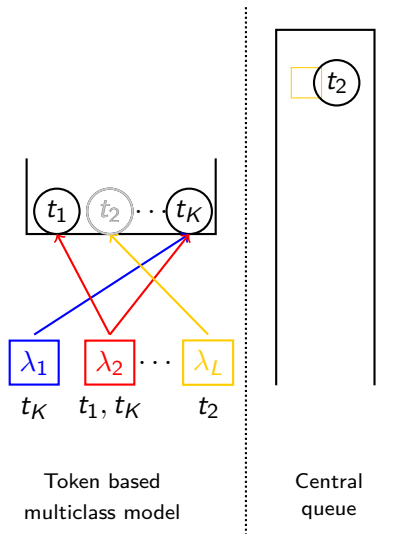
- A single central queue
- A set $\mathcal{T} = \{t_1, \dots, t_K\}$ of K tokens
- Only jobs with tokens are served
- There are L job classes
- Poisson arrival rate λ_i for class i
- A compatibility graph (arbitrary)
- Features of token based queue
 - Token assignment
 - Releasing a token
 - Token service rate
 - State space



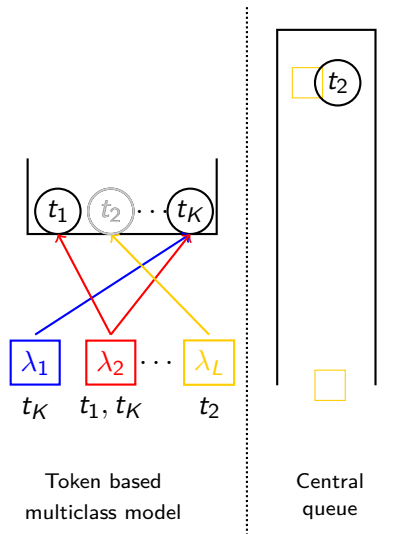
- An arriving job must pick a compatible token if available



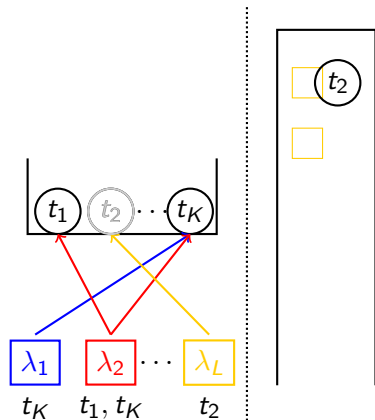
- An arriving job must pick a compatible token if available



- An arriving job must pick a compatible token if available



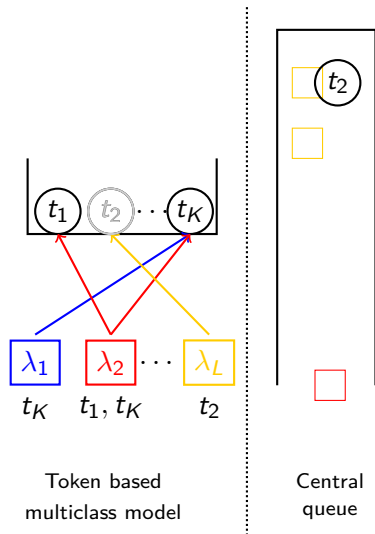
- An arriving job must pick a compatible token if available



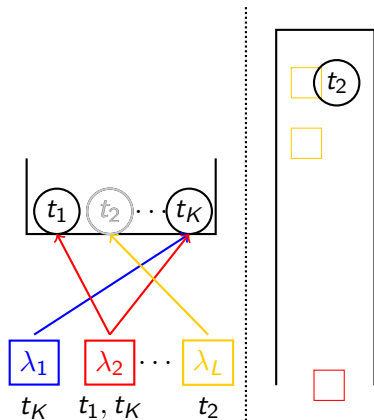
Token based
multiclass model

Central
queue

- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue



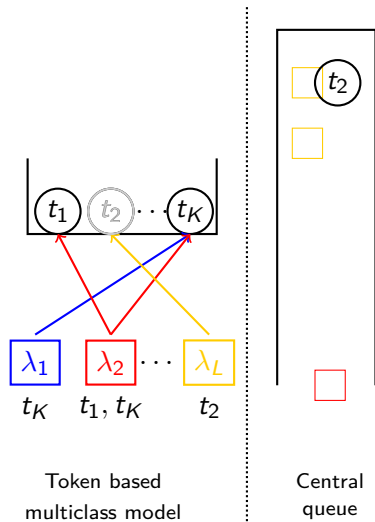
- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue



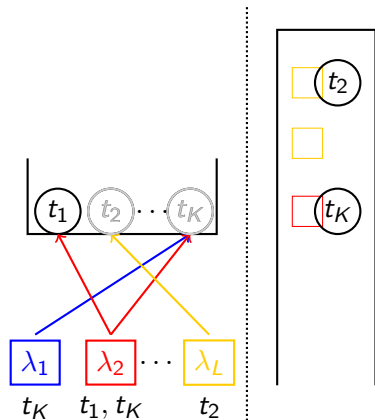
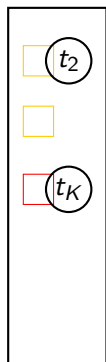
Token based
multiclass model

Central
queue

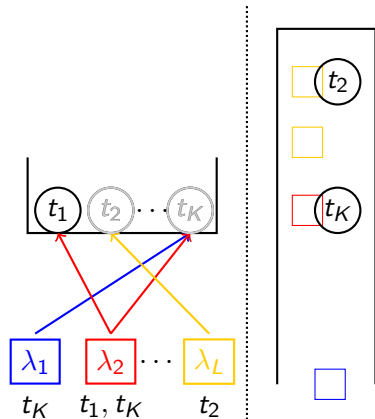
- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token



- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token
- An *assignment rule* specifies the tie-breaking rule

Token based
multiclass modelCentral
queue

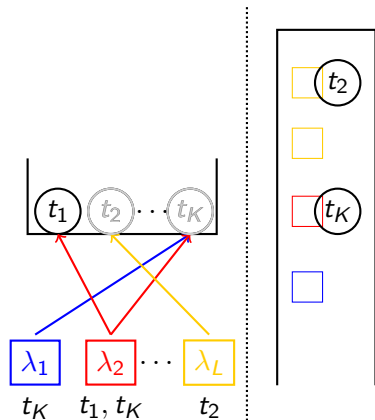
- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token
- An *assignment rule* specifies the tie-breaking rule



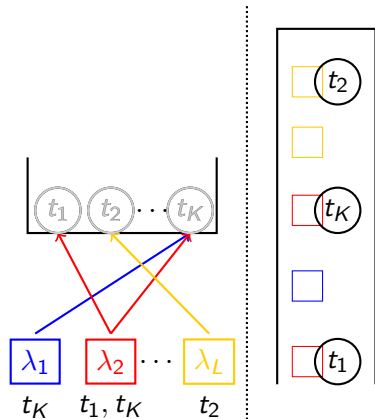
Token based
multiclass model

Central
queue

- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token
- An *assignment rule* specifies the tie-breaking rule

Token based
multiclass modelCentral
queue

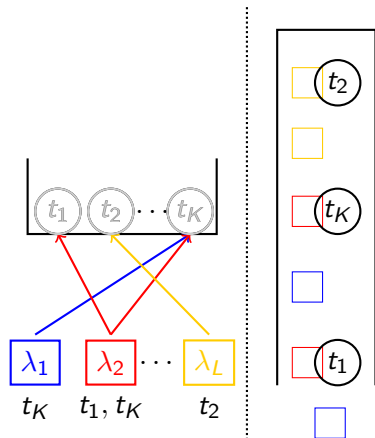
- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token
- An *assignment rule* specifies the tie-breaking rule



Token based
multiclass model

Central
queue

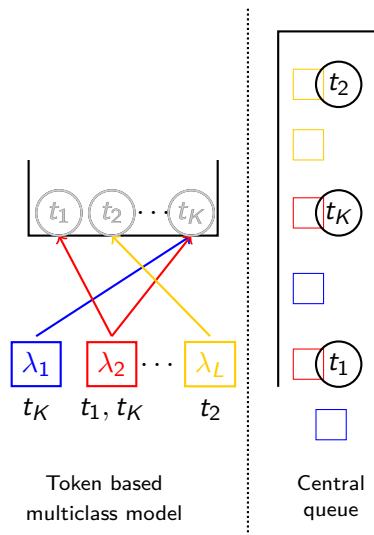
- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token
- An *assignment rule* specifies the tie-breaking rule



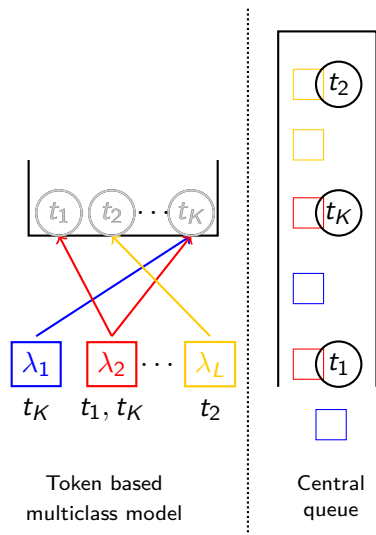
Token based
multiclass model

Central
queue

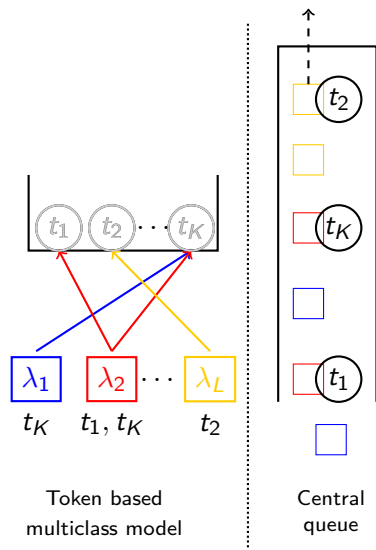
- An arriving job must pick a compatible token if available
- Job with no feasible token will wait in the queue
- A job can claim only one compatible token
- An *assignment rule* specifies the tie-breaking rule



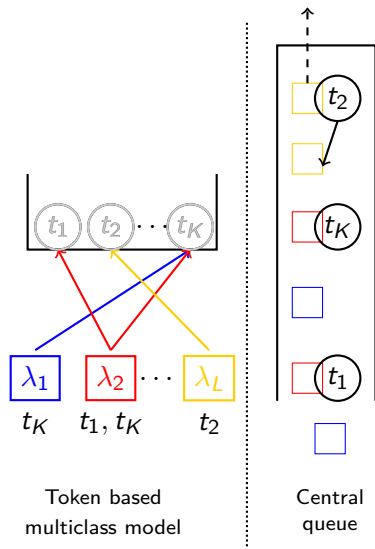
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate



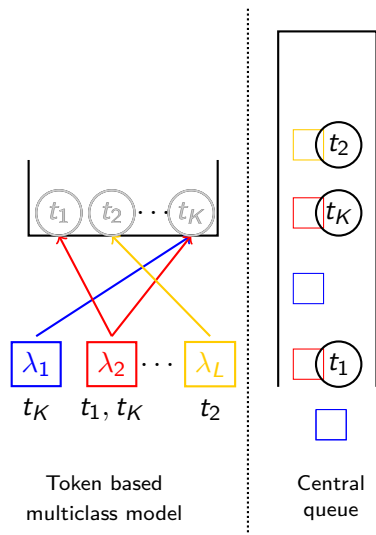
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job



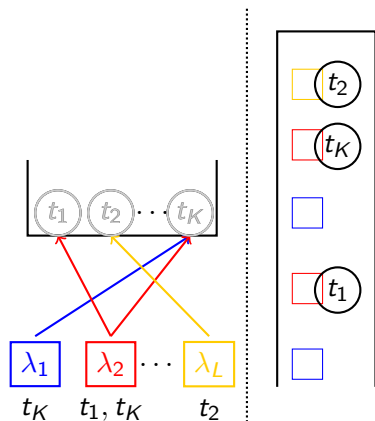
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job



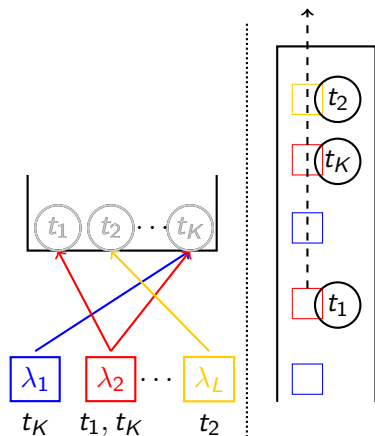
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job



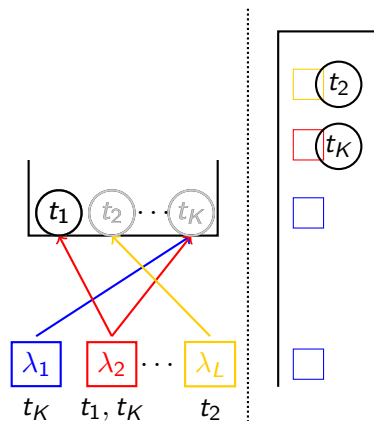
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job

Token based
multiclass modelCentral
queue

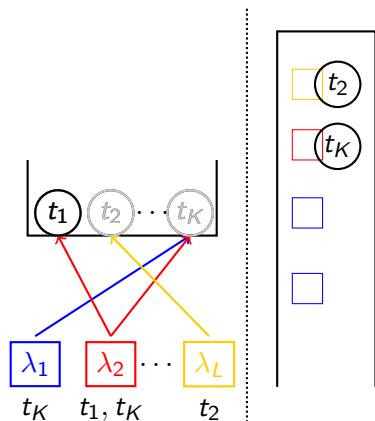
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job

Token based
multiclass modelCentral
queue

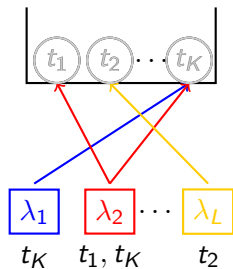
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job

Token based
multiclass modelCentral
queue

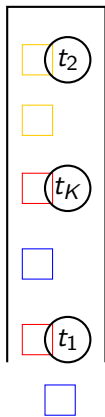
- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job
- If no waiting compatible job present, token added back to the token set

Token based
multiclass modelCentral
queue

- Jobs have an exponential service requirement with unit mean
- Only jobs with tokens will receive a non-negative service rate
- When a job departs, the released token picks the next compatible waiting job
- If no waiting compatible job present, token added back to the token set

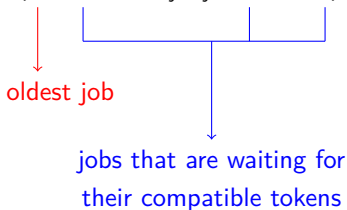


Token based
multiclass model



Central
queue

- Markovian descriptor for token based model
- Anonymize the class information
- $(t_2, 1, t_k, 1, t_1, 1)$
- $(T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$

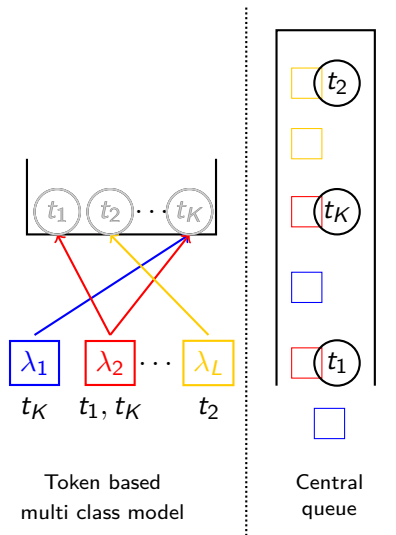


Token assignment

Releasing a token

Token service rate

State space

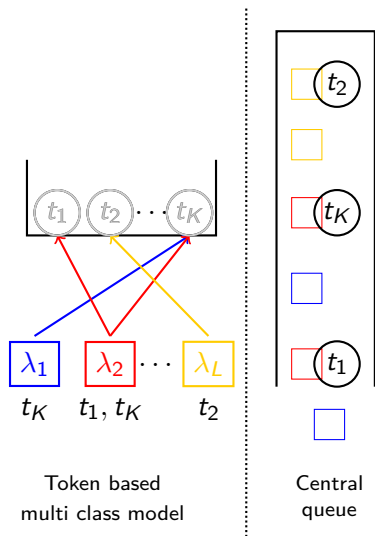


Token assignment

Releasing a token

Token service rate

State space



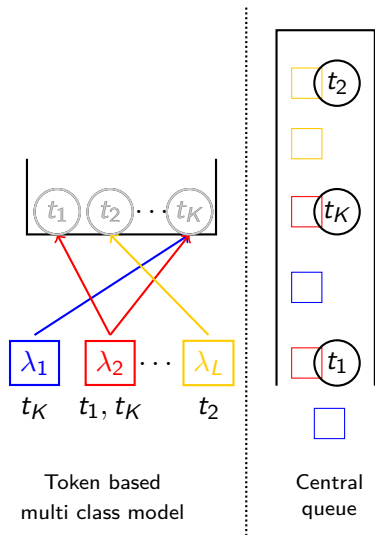
- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$

Token assignment

Releasing a token

Token service rate

State space



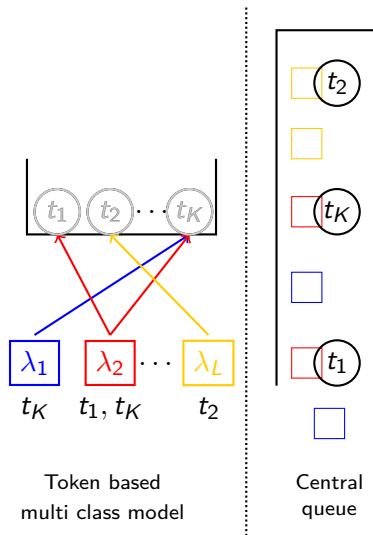
- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- $\mu_{T_j}(x)$ denote the departure rate of the job with token T_j

Token assignment

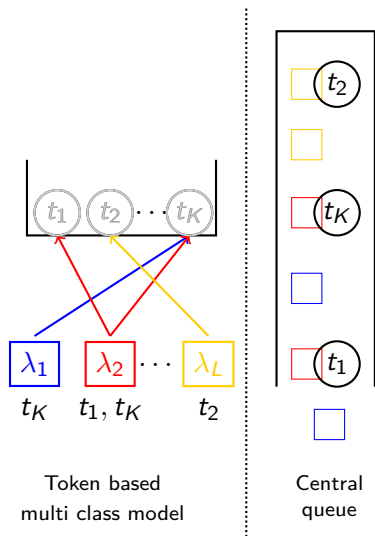
Releasing a token

Token service rate

State space



- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- $\mu_{T_j}(x)$ denote the departure rate of the job with token T_j
- $\mu(x) = \sum_{j=1}^i \mu_{T_j}(x)$



- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- $\mu_{T_j}(x)$ denote the departure rate of the job with token T_j
- $\mu(x) = \sum_{j=1}^i \mu_{T_j}(x)$

Token service rate

Objects of departures T_j

$$s = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$$

¹²A. Krzesinski, Order independent queues, in "Queueing Networks: a fundamental approach", Eds: R. Boucherie, N. van Dijk, 2011

Token service rate

Objects of departures T_j

$$s = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$$

Require the departure function $\hat{\mu}(\cdot)$ satisfying the following

1. $\hat{\mu}(T_1 \dots T_j) - \hat{\mu}(T_1 \dots T_{j-1}) \geq 0$
2. $\hat{\mu}(T_1 \dots T_j) = \hat{\mu}(T_{\sigma(1)} \dots T_{\sigma(j-1)})$

Order Independent rates¹² \implies Sufficient condition for product-form

¹²A. Krzesinski, Order independent queues, in "Queueing Networks: a fundamental approach", Eds: R. Boucherie, N. van Dijk, 2011

Token service rate

Objects of departures T_j

$$s = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$$

Require the departure function $\hat{\mu}(\cdot)$ satisfying the following

1. $\hat{\mu}(T_1 \dots T_j) - \hat{\mu}(T_1 \dots T_{j-1}) \geq 0$
2. $\hat{\mu}(T_1 \dots T_j) = \hat{\mu}(T_{\sigma(1)} \dots O_{\sigma(j-1)})$

Order Independent rates¹² \implies Sufficient condition for product-form

For redundancy-*d* c.o.s., the T_j correspond to server M_j
For redundancy-*d* c.o.c., the T_j correspond to first customer of a class c_j

¹²A. Krzesinski, Order independent queues, in "Queueing Networks: a fundamental approach", Eds: R. Boucherie, N. van Dijk, 2011

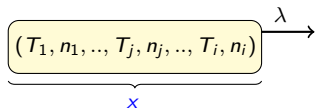
State transitions

$$(T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$$

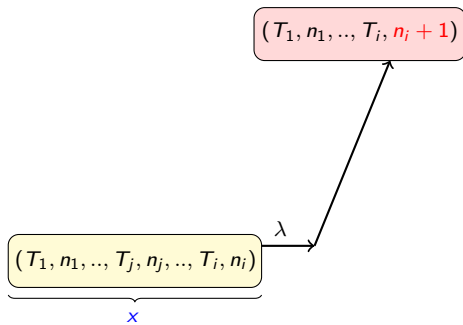
State transitions

$$\underbrace{(T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)}_x$$

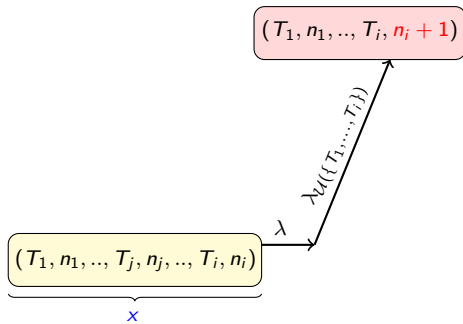
State transitions



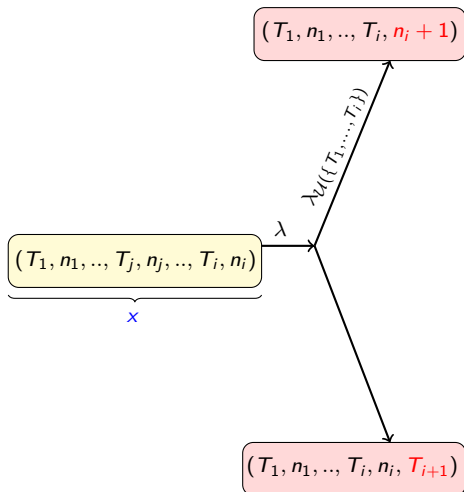
State transitions



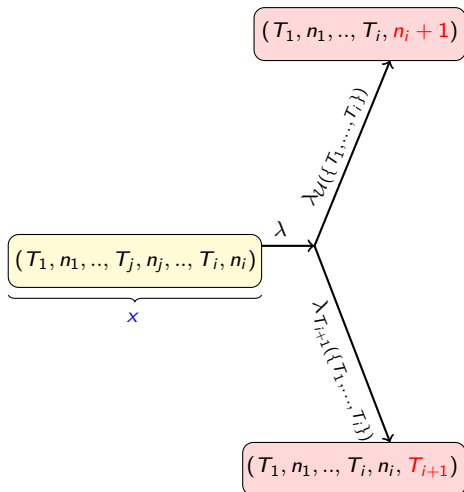
State transitions



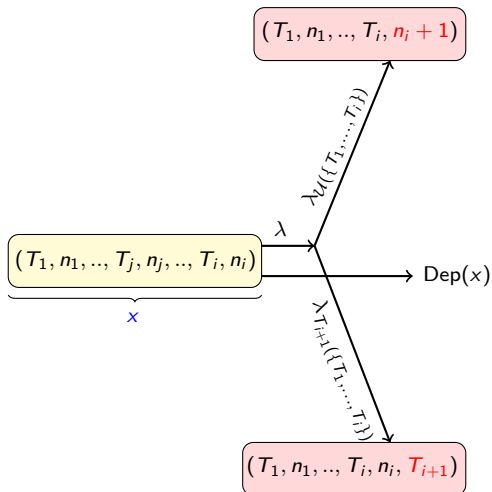
State transitions



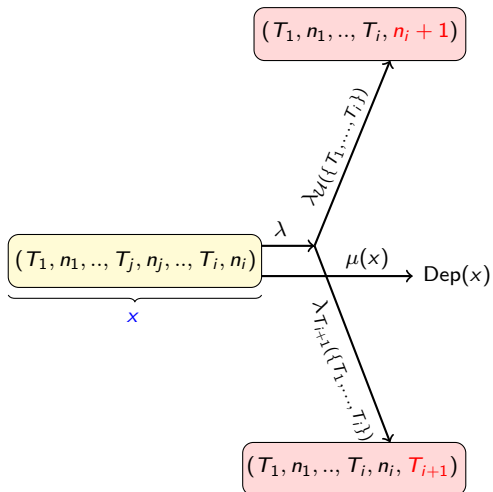
State transitions



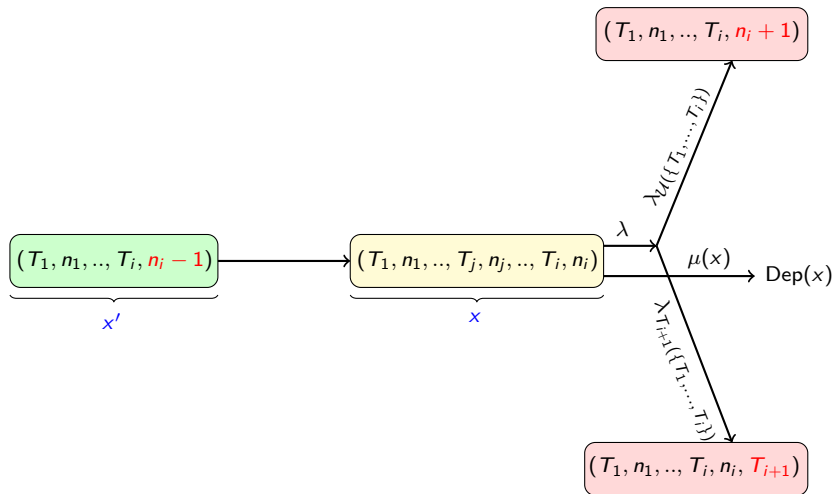
State transitions



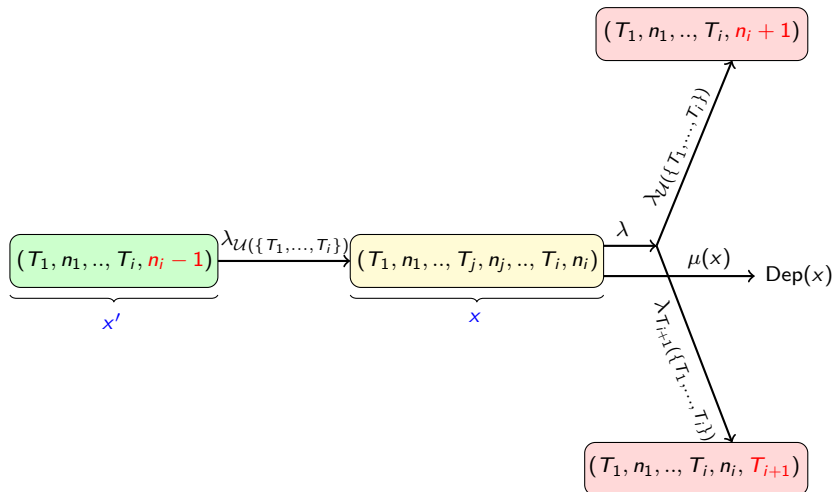
State transitions



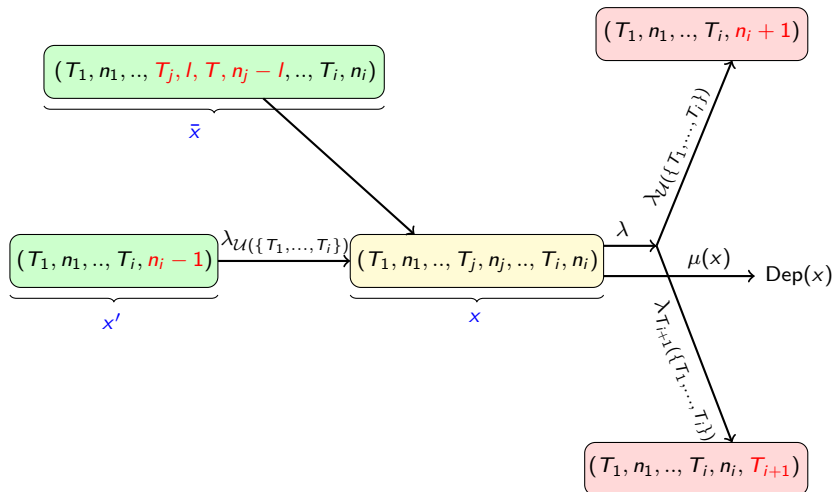
State transitions



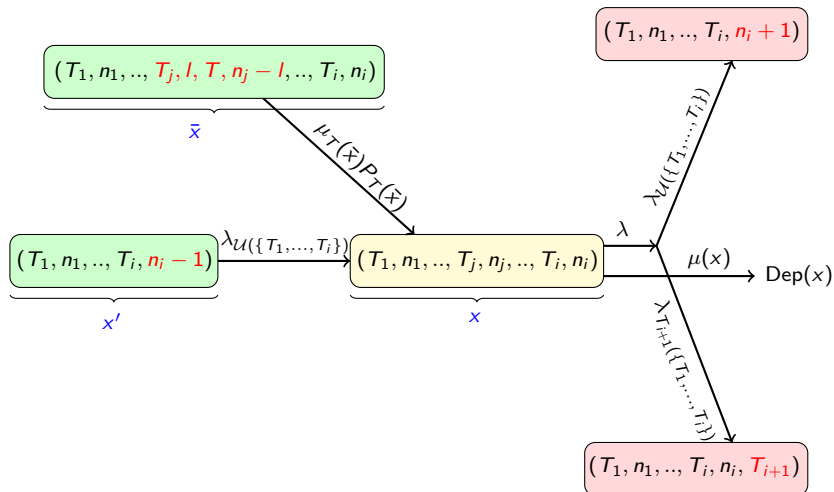
State transitions



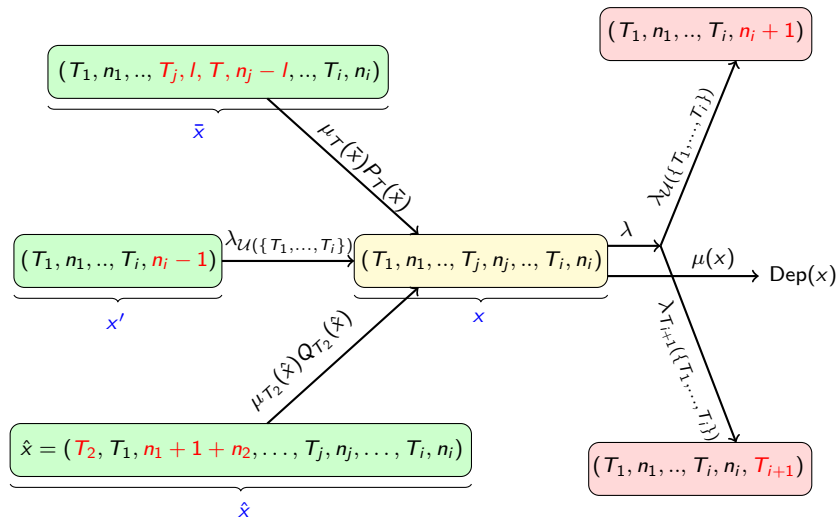
State transitions



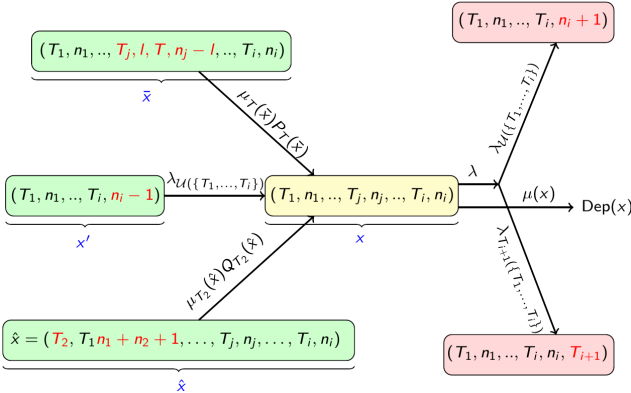
State transitions



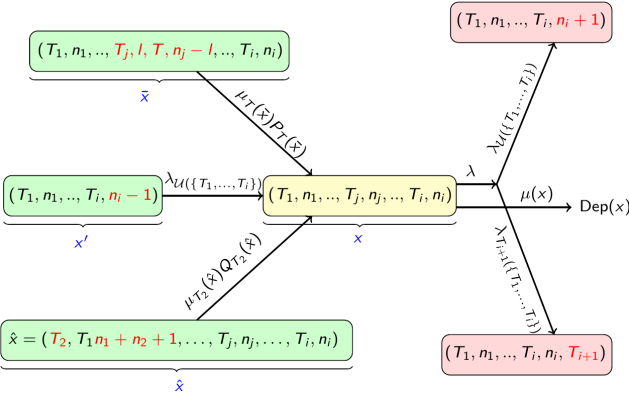
State transitions



Global balance equations (GBE) for stationary distribution

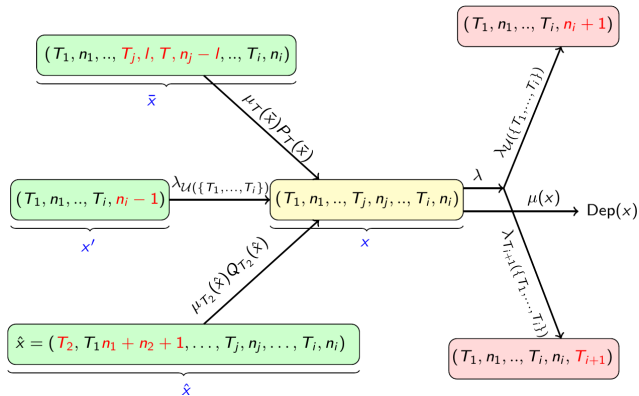


Global balance equations (GBE) for stationary distribution



- **GBE:** Average rate of leaving state x = Average rate of entering state x

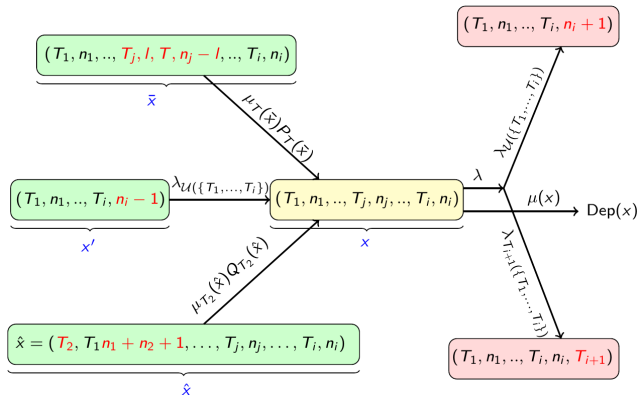
Global balance equations (GBE) for stationary distribution



- **GBE:** Average rate of leaving state x = Average rate of entering state x

$$\pi(x) \left[\lambda_{\mathcal{U}(\{\cdot\})} + \sum_{T_{i+1}} \lambda_{T_{i+1}}(\cdot) + \mu(x) \right]$$

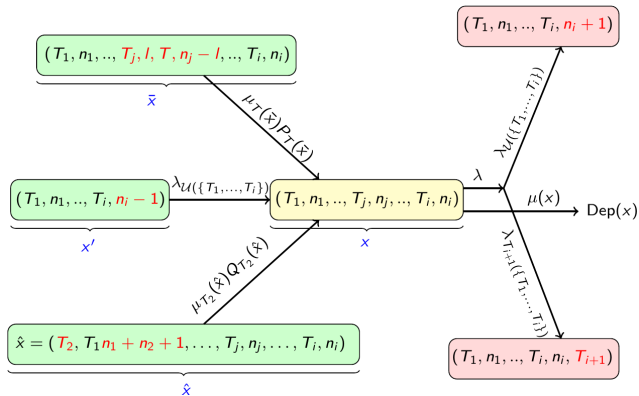
Global balance equations (GBE) for stationary distribution



- **GBE:** Average rate of leaving state x = Average rate of entering state x

$$\pi(x) \left[\lambda_{\mathcal{U}(\{\cdot\})} + \sum_{T_{i+1}} \lambda_{T_{i+1}}(\cdot) + \mu(x) \right] = \pi(x') \lambda_{\mathcal{U}(\cdot)} + \sum \pi(\bar{x}) \mu_T(\bar{x}) P_T(\bar{x}) + \sum \pi(\hat{x}) \mu_{T_2}(\hat{x}) Q_{T_2}(\hat{x})$$

Global balance equations (GBE) for stationary distribution



- **GBE:** Average rate of leaving state x = Average rate of entering state x

$$\pi(x) \left[\lambda_{\mathcal{U}(\{\cdot\})} + \sum_{T_{i+1}} \lambda_{T_{i+1}}(\cdot) + \mu(x) \right] = \pi(x') \lambda_{\mathcal{U}(\cdot)} + \sum \pi(\bar{x}) \mu_T(\bar{x}) P_T(\bar{x}) + \sum \pi(\hat{x}) \mu_{T_2}(\hat{x}) Q_{T_2}(\hat{x})$$

- We solve this using **partial balance equations (PBE's)**.

Main Result

Theorem

The steady state distribution for the token based model in state

$x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

Main Result

Theorem

The steady state distribution for the token based model in state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

$$\pi(x) = \pi(0) \prod_{j=1}^i \left(\frac{\lambda_{T_j}(\{T_1, \dots, T_{j-1}\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right) \prod_{j=1}^i \left(\frac{\lambda_{\mathcal{U}(\{T_1, \dots, T_j\})}}{\sum_{l=1}^j \mu_{T_l}(x)} \right)^{n_j}$$

Main Result

Theorem

The steady state distribution for the token based model in state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

$$\pi(x) = \pi(0) \prod_{j=1}^i \left(\frac{\lambda_{T_j}(\{T_1, \dots, T_{j-1}\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right) \prod_{j=1}^i \left(\frac{\lambda_{U(\{T_1, \dots, T_j\})}}{\sum_{l=1}^j \mu_{T_l}(x)} \right)^{n_j}$$

↓
normalising
constant

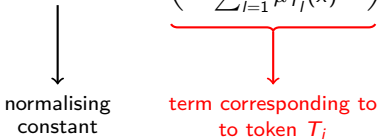
Main Result

Theorem

The steady state distribution for the token based model in state

$x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

$$\pi(x) = \pi(0) \prod_{j=1}^i \left(\frac{\lambda_{T_j}(\{T_1, \dots, T_{j-1}\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right) \prod_{j=1}^i \left(\frac{\lambda_{\mathcal{U}(\{T_1, \dots, T_j\})}}{\sum_{l=1}^j \mu_{T_l}(x)} \right)^{n_j}$$



normalising constant

term corresponding to token T_j

Main Result

Theorem

The steady state distribution for the token based model in state

$x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

$$\pi(x) = \pi(0) \prod_{j=1}^i \left(\frac{\lambda_{T_j}(\{T_1, \dots, T_{j-1}\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right) \prod_{j=1}^i \left(\frac{\lambda_{U}(\{T_1, \dots, T_j\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right)^{n_j}$$

normalising constant

term corresponding to token T_j

term corresponding to each of the n_j waiting jobs

Main Result

Theorem

The steady state distribution for the token based model in state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

$$\pi(x) = \pi(0) \prod_{j=1}^i \left(\frac{\lambda_{T_j}(\{T_1, \dots, T_{j-1}\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right) \prod_{j=1}^i \left(\frac{\lambda_{U(\{T_1, \dots, T_j\})}}{\sum_{l=1}^j \mu_{T_l}(x)} \right)^{n_j}$$

normalising constant term corresponding to token T_j term corresponding to each of the n_j waiting jobs

- PGF for the number of waiting jobs and the total number of jobs

Main Result

Theorem

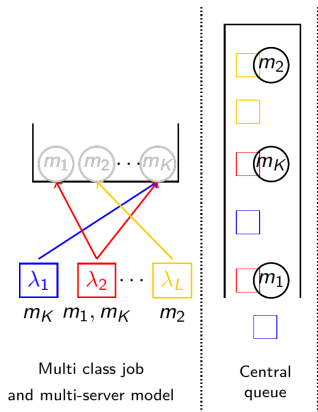
The steady state distribution for the token based model in state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$ is given by

$$\pi(x) = \pi(0) \prod_{j=1}^i \left(\frac{\lambda_{T_j}(\{T_1, \dots, T_{j-1}\})}{\sum_{l=1}^j \mu_{T_l}(x)} \right) \prod_{j=1}^i \left(\frac{\lambda_{U(\{T_1, \dots, T_j\})}}{\sum_{l=1}^j \mu_{T_l}(x)} \right)^{n_j}$$

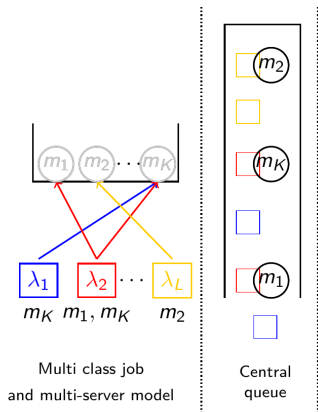
normalising constant term corresponding to token T_j term corresponding to each of the n_j waiting jobs

- PGF for the number of waiting jobs and the total number of jobs
- PGF for the waiting time and sojourn time for a job

Visschers and order-independent (OI) models

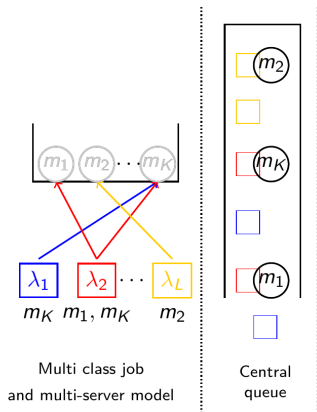


Visschers and order-independent (OI) models



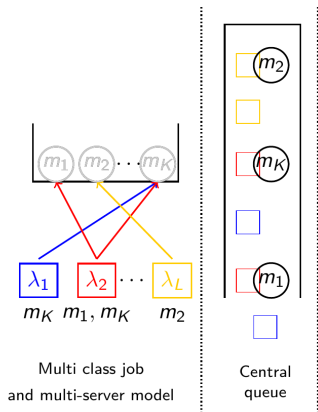
- Tokens replaced by servers

Visschers and order-independent (OI) models



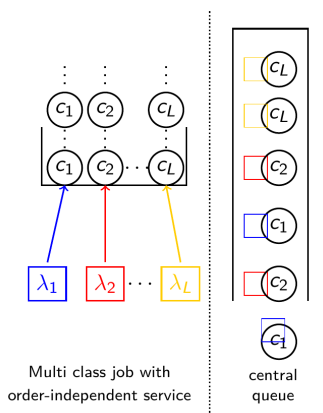
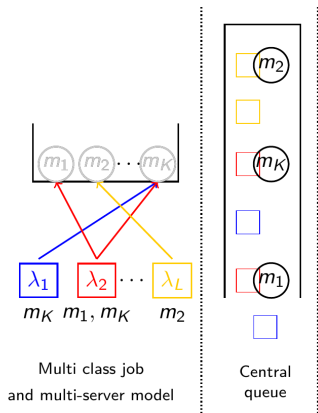
- Tokens replaced by servers
- Constant service rate

Visschers and order-independent (OI) models



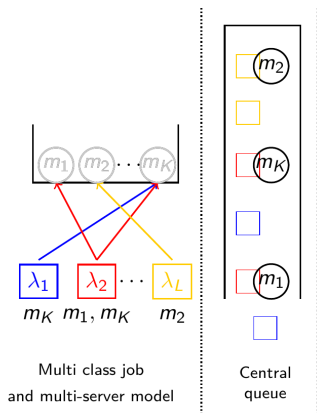
- Tokens replaced by servers
- Constant service rate
- *Assignment rule*

Visschers and order-independent (OI) models

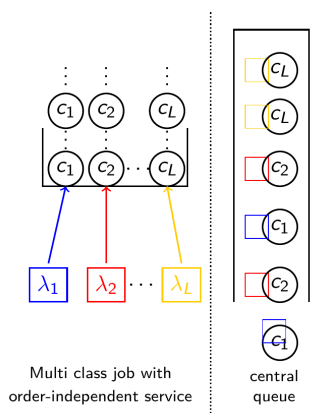


- Tokens replaced by servers
- Constant service rate
- *Assignment rule*

Visschers and order-independent (OI) models

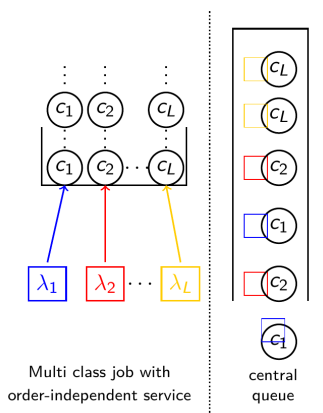
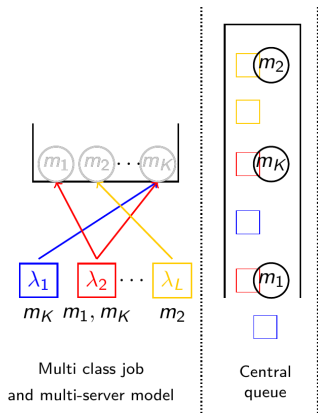


- Tokens replaced by servers
- Constant service rate
- *Assignment rule*



- Each job of class i gets a token c_i

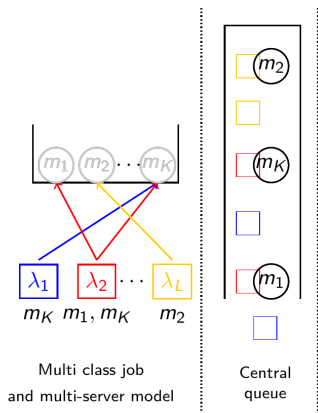
Visschers and order-independent (OI) models



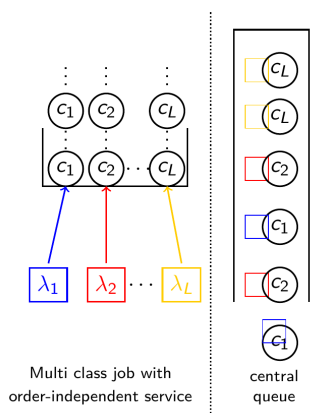
- Tokens replaced by servers
- Constant service rate
- *Assignment rule*

- Each job of class i gets a token c_i
- Order-independent (OI) service rate

Visschers and order-independent (OI) models

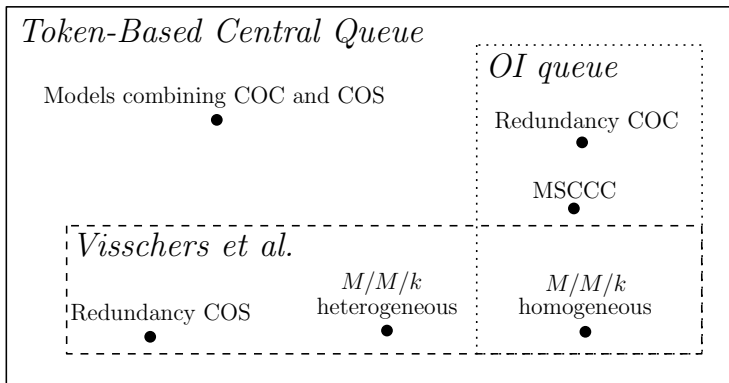


- Tokens replaced by servers
- Constant service rate
- *Assignment rule*

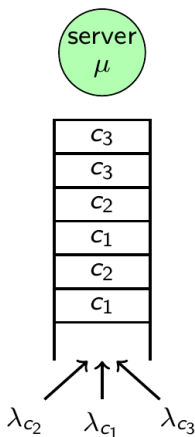


- Each job of class i gets a token c_i
- Order-independent (OI) service rate
- No *assignment rule*

A classification of token-based central queues

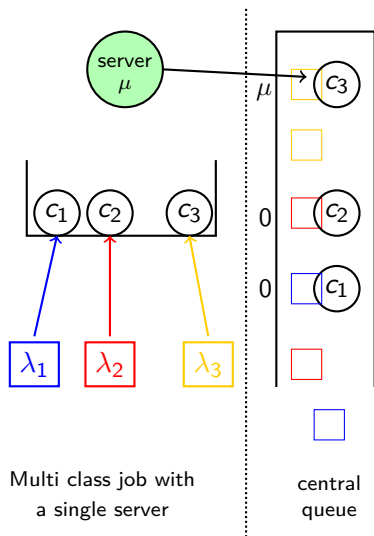


Example 1: Multiclass $M/M/1$ queue



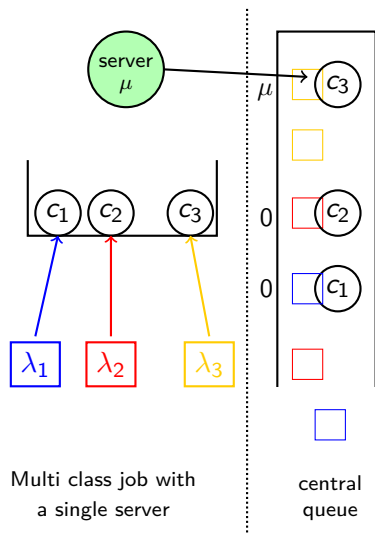
- single server with multiple classes
- exponential service requirement
- state space of the form $(c_3, c_3, c_2, c_1, c_2, c_1)$
- $\pi(c_3, c_3, c_2, c_1, c_2, c_1)$
 $= (1 - \rho) \rho_{c_3}^2 \rho_{c_2} \rho_{c_1} \rho_{c_2} \rho_{c_1}$
where $\rho_i = \frac{\lambda_i}{\mu}$

Example 1: Multiclass $M/M/1$ queue (token view)



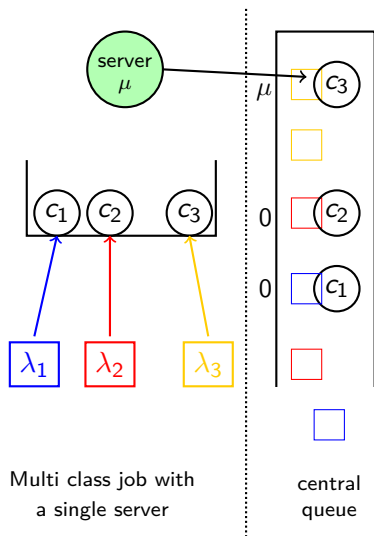
- We associate a unique token per class
- state space $(c_3, 1, c_2, c_1, 2)$
- Service rate μ for the first token
service rate 0 for the rest

Example 1: Multiclass $M/M/1$ queue (token view)



- We associate a unique token per class
- state space $(c_3, 1, c_2, c_1, 2)$
- Service rate μ for the first token
service rate 0 for the rest

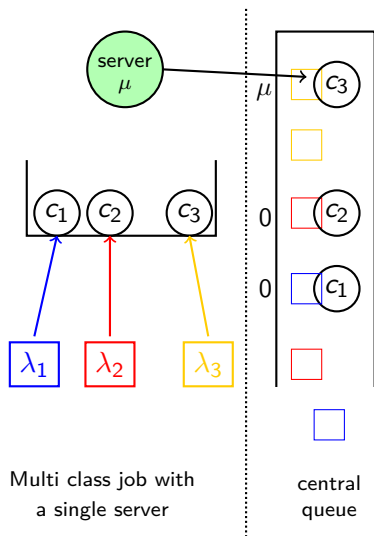
Example 1: Multiclass $M/M/1$ queue (token view)



- We associate a unique token per class
- state space $(c_3, 1, c_2, c_1, 2)$
- Service rate μ for the first token
service rate 0 for the rest

- $\pi(c_3, 1, c_2, c_1, 2) =$
 $(1 - \rho) \left(\frac{\lambda_{c_3}}{\mu}\right)^2 \left(\frac{\lambda_{c_2}}{\mu}\right) \left(\frac{\lambda_{c_1}}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^2$

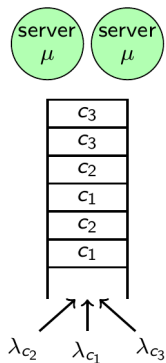
Example 1: Multiclass $M/M/1$ queue (token view)



- We associate a unique token per class
- state space $(c_3, 1, c_2, c_1, 2)$
- Service rate μ for the first token
service rate 0 for the rest

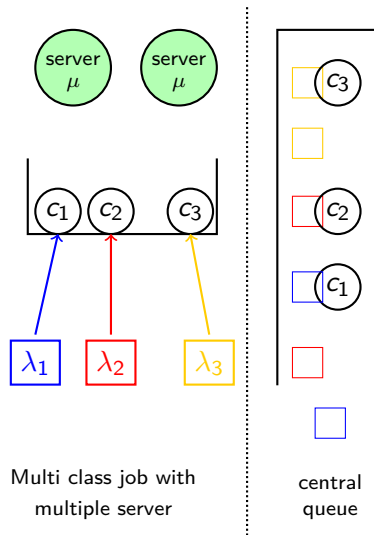
- $\pi(c_3, 1, c_2, c_1, 2) =$
 $(1 - \rho) \left(\frac{\lambda_{c_3}}{\mu}\right)^2 \left(\frac{\lambda_{c_2}}{\mu}\right) \left(\frac{\lambda_{c_1}}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^2$

Example 2: Multi-server station with concurrent control of customers (MSCCC)



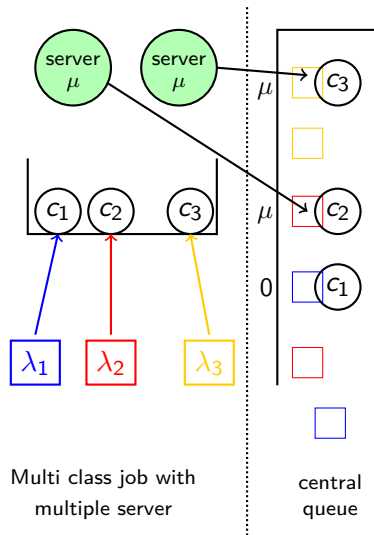
- 2 identical servers and 3 classes
- Servers are compatible with all classes
- Both servers cannot serve the same class at any time (**concurrent control**)

Example 2: Multi-server station with concurrent control of customers (MSCCC)



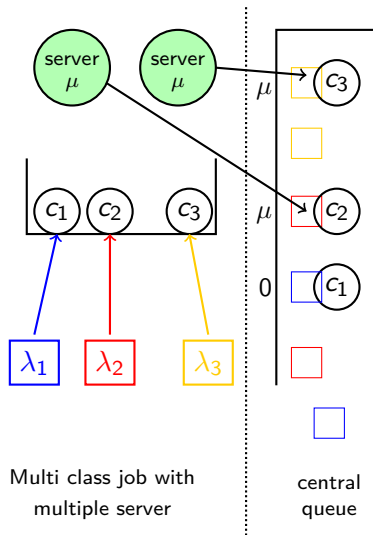
- One token per class as earlier
- At most 1 job per class is served (**concurrent control**)
- Two server, three tokens
- Associate the two servers with the first 2 active tokens

Example 2: Multi-server station with concurrent control of customers (MSCCC)



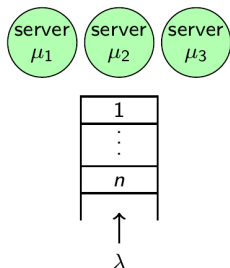
- One token per class as earlier
- At most 1 job per class is served (**concurrent control**)
- Two server, three tokens
- Associate the two servers with the first 2 active tokens
- Order-independent service rate where only the first 2 token have a service rate μ

Example 2: Multi-server station with concurrent control of customers (MSCCC)



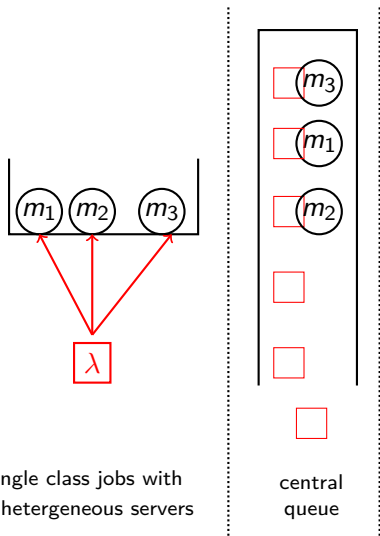
- One token per class as earlier
- At most 1 job per class is served (**concurrent control**)
- Two server, three tokens
- Associate the two servers with the first 2 active tokens
- Order-independent service rate where only the first 2 token have a service rate μ
- $\pi(c_3, 1, c_2, c_1, 2) = (1 - \rho) \left(\frac{\lambda_{c_3}}{\mu}\right)^2 \left(\frac{\lambda_{c_2}}{2\mu}\right) \left(\frac{\lambda_{c_1}}{2\mu}\right) \left(\frac{\lambda}{2\mu}\right)^2$

Example 3: $M/M/K$ queue

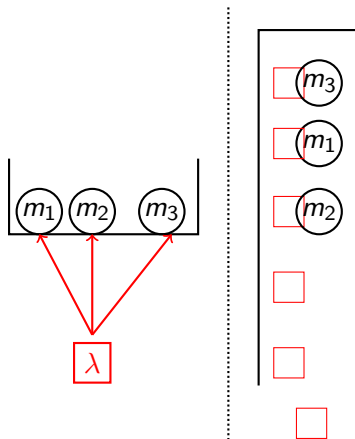


- One class of jobs and K possibly heterogeneous servers
- At most 3 jobs can be served at a time
- Also known as *Erlang-K* model
- When servers are identical, state space denoted by number of jobs in the system
- $\pi(n) = \pi(0) \frac{\rho^n K^K}{K!}$ where $\rho = \frac{\lambda}{K\mu}$

Example 3: $M/M/K$ queue



Example 3: $M/M/K$ queue

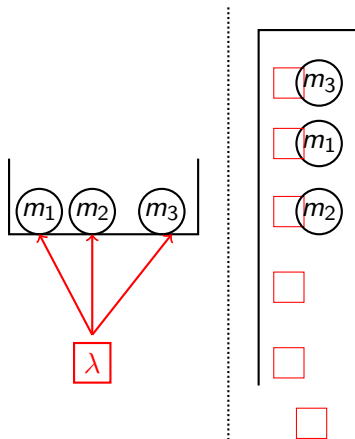


single class jobs with
 K heterogeneous servers

central
queue

- Single class & token m_i for server i

Example 3: $M/M/K$ queue

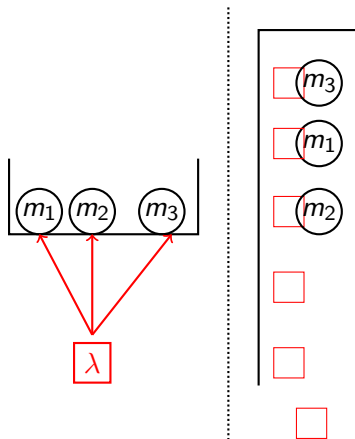


single class jobs with
 K heterogeneous servers

central
queue

- Single class & token m_i for server i
- Arriving jobs are compatible with all tokens

Example 3: $M/M/K$ queue

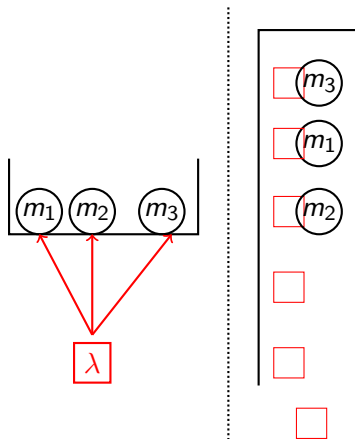


single class jobs with
 K heterogeneous servers

central
queue

- Single class & token m_i for server i
- Arriving jobs are compatible with all tokens
- *Assignment rule is uniform*

Example 3: $M/M/K$ queue

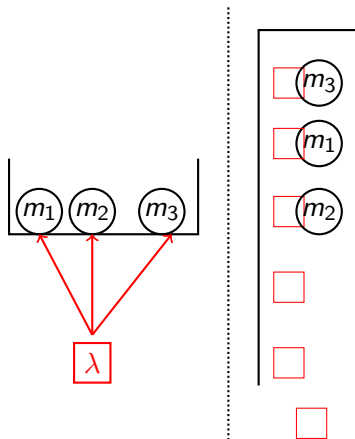


single class jobs with
 K heterogeneous servers

central
queue

- Single class & token m_i for server i
- Arriving jobs are compatible with all tokens
- *Assignment rule is uniform*
- $(m_3, m_1, m_2, 3)$

Example 3: $M/M/K$ queue



single class jobs with
 K heterogeneous servers

central
queue

- Single class & token m_i for server i
- Arriving jobs are compatible with all tokens
- Assignment rule is uniform

- $(m_3, m_1, m_2, 3)$

- $\pi(m_3, m_1, m_2, 3) =$

$$\pi(0) \frac{\lambda}{3\mu_3} \frac{\lambda}{2(\mu_3+\mu_1)} \frac{\lambda}{(\mu_3+\mu_1+\mu_2)} \left(\frac{\lambda}{(\mu_3+\mu_1+\mu_2)} \right)^3$$

Outline

- ▶ Redundancy
- ▶ Central Queue Architecture
- ▶ Order Independent Descriptor
 - ▶ Redundancy and cancel on complete
- ▶ Aggregated State Descriptor
 - ▶ Redundancy and cancel on start
- ▶ Generalizations:
 - ▶ Token-based framework
 - ▶ Generalized Order Independent
- ▶ Impact of assumptions: scheduling and independence

Generalized OI queue¹³

Detailed state descriptor \vec{z}_m , where m is number of jobs

- ▶ It includes jobs in **service** and in the **queue**

service \implies track server

queue \implies track class

¹³K. Gardner, R. Righter, Product Forms for FCFS Queueing Models with Arbitrary Server-Job Compatibilities: An Overview, to appear in QUESTA

Generalized OI queue¹³

Detailed state descriptor \vec{z}_m , where m is number of jobs

- ▶ It includes jobs in **service** and in the **queue**

service \implies track server

queue \implies track class

- ▶ $\mu(\vec{z}_m)$ satisfies the OI Properties

Theorem: The steady-state distribution of the Generalized OI queue is product form

\implies Unifying framework for product-form distributions

¹³K. Gardner, R. Righter, Product Forms for FCFS Queueing Models with Arbitrary Server-Job Compatibilities: An Overview, to appear in QUESTA

Outline

- ▶ Redundancy
- ▶ Central Queue Architecture
- ▶ Order Independent Descriptor
 - ▶ Redundancy and cancel on complete
- ▶ Aggregated State Descriptor
 - ▶ Redundancy and cancel on start
- ▶ Generalizations:
 - ▶ Token-based framework
 - ▶ Generalized Order Independent
- ▶ Impact of assumptions: scheduling and independence

Stability: Impact of independence assumption

- ▶ With i.i.d. copies and FCFS, redundancy does not impact stability.

⇒ papers by Gardner et al., Bonald et al.

$d = K \implies$ *single server with rate μK*

$d = 1 \implies$ *K indep. single servers with rate μ*

Stability: Impact of independence assumption

- ▶ With i.i.d. copies and FCFS, redundancy does not impact stability.

⇒ papers by Gardner et al., Bonald et al.

$d = K \implies$ *single server with rate μK*

$d = 1 \implies$ *K indep. single servers with rate μ*

- ▶ What without i.i.d. assumption?

$d = K \implies$ *single server with rate μ*

$d = 1 \implies$ *K indep. single servers with rate μ*

Stability of redundancy: Impact of assumptions

Most of existing literature is with i.i.d. copies and FCFS

⇒ stability not reduced

yield results that are qualitatively misleading.

Need of better models: *S&X*¹⁴

¹⁴Gardner, Harchol-Balter, Scheller-Wolf, Van Houdt, A Better Model for Job Redundancy: Decoupling Server Slowdown and Job Size, IEEE/ACM ToN, 2017

Stability of redundancy: Impact of assumptions

Most of existing literature is with i.i.d. copies and FCFS

⇒ stability not reduced

yield results that are qualitatively misleading.

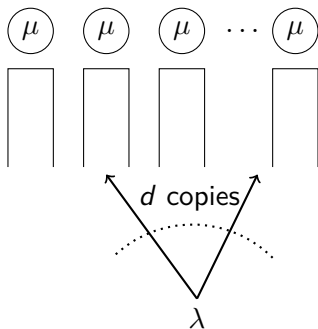
Need of better models: $S&X$ ¹⁴

Main questions:

- ▶ How does redundancy with identical copies impact stability?
- ▶ Does stability depend on the scheduling discipline?

¹⁴Gardner, Harchol-Balter, Scheller-Wolf, Van Houdt, A Better Model for Job Redundancy: Decoupling Server Slowdown and Job Size, IEEE/ACM ToN, 2017

Redundancy d

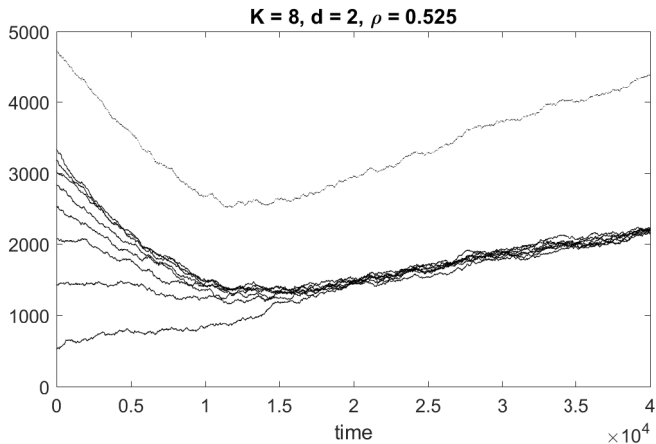


We study¹⁵:

- ▶ c.o.c. with identical copies
- ▶ scheduling discipline in servers can be PS, FCFS, or ROS.

¹⁵E. Anton, U. Ayesta, M. Jonckheere, I.M. Verloop. On the stability of redundancy models, to appear in OR

Instability with identical copies



Neither total number nor minimum are Lyapunov functions

Scheduling disciplines

- ▶ FCFS.
- ▶ Processor-Sharing (PS). The capacity is shared equally among all copies
- ▶ Random Order of Service (ROS)

¹⁶ Baccelli, Foss, On the Saturation Rule for the Stability of Queues, JAP, 1995

Scheduling disciplines

- ▶ FCFS.
- ▶ Processor-Sharing (PS). The capacity is shared equally among all copies
- ▶ Random Order of Service (ROS)

FCFS. Stability through saturated system¹⁶

¹⁶ Baccelli, Foss, On the Saturation Rule for the Stability of Queues, JAP, 1995

Scheduling disciplines

- ▶ FCFS.
- ▶ Processor-Sharing (PS). The capacity is shared equally among all copies
- ▶ Random Order of Service (ROS)

FCFS. Stability through saturated system¹⁶

For PS.

Lower Bound: Every instant, put the copy with highest attained service, in the server with smallest number of copies.

Upper Bound: All copies need to be served for the job to be completed.

¹⁶ Baccelli, Foss, On the Saturation Rule for the Stability of Queues, JAP, 1995

Summary of stability conditions

	PS	FCFS	ROS	Priority policy
i.i.d	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda \ll \mu K$
i.c.	$\lambda < \mu \frac{K}{d}$	$\lambda < \bar{\ell} \mu$	$\lambda < \mu K$	-

¹⁷Y. Raaijmakers, S. Borst, O. Boxma: Delta probing policies for redundancy. Perform. Eval (2018)

Summary of stability conditions

	PS	FCFS	ROS	Priority policy
i.i.d	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda \ll \mu K$
i.c.	$\lambda < \mu \frac{K}{d}$	$\lambda < \bar{\ell} \mu$	$\lambda < \mu K$	-

- ▶ With PS, stability condition is the same as if all copies had to be served.
- ▶ The stability region is larger for FCFS than for PS.
- ▶ ROS. No stability reduction.

¹⁷Y. Raaijmakers, S. Borst, O. Boxma: Delta probing policies for redundancy. Perform. Eval (2018)

Summary of stability conditions

	PS	FCFS	ROS	Priority policy
i.i.d	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda \ll \mu K$
i.c.	$\lambda < \mu \frac{K}{d}$	$\lambda < \bar{\ell} \mu$	$\lambda < \mu K$	-

- ▶ With PS, stability condition is the same as if all copies had to be served.
- ▶ The stability region is larger for FCFS than for PS.
- ▶ ROS. No stability reduction.
- ▶ Stability depends on scheduling discipline.

¹⁷Y. Raaijmakers, S. Borst, O. Boxma: Delta probing policies for redundancy. Perform. Eval (2018)

Summary of stability conditions

	PS	FCFS	ROS	Priority policy
i.i.d	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda < \mu K$	$\lambda \ll \mu K$
i.c.	$\lambda < \mu \frac{K}{d}$	$\lambda < \bar{\ell} \mu$	$\lambda < \mu K$	-

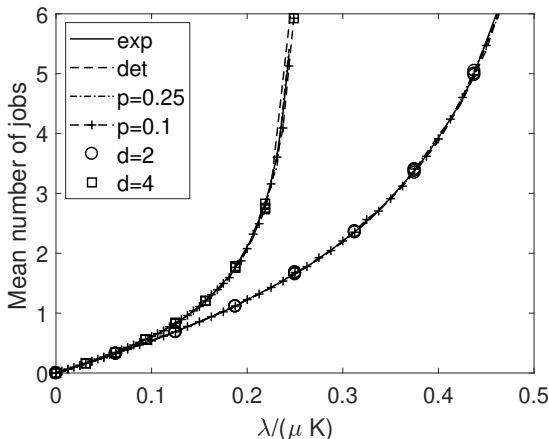
- ▶ With PS, stability condition is the same as if all copies had to be served.
- ▶ The stability region is larger for FCFS than for PS.
- ▶ ROS. No stability reduction.
- ▶ Stability depends on scheduling discipline.

Need to develop strategies that preserve stability¹⁷

¹⁷Y. Raaijmakers, S. Borst, O. Boxma: Delta probing policies for redundancy. Perform. Eval (2018)

Non-exponential service requirements: PS¹⁸

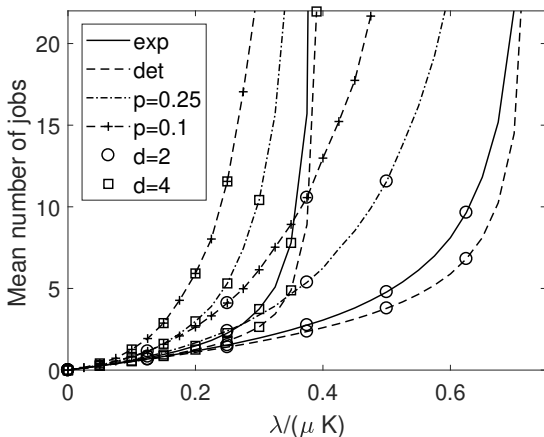
Mean number of jobs with identical copies and exponential, deterministic and degenerate hyperexponential service requirements.



¹⁸PG Taylor, Insensitivity in Stochastic Models, in "Queueing Networks: a fundamental approach", Eds: R. Boucherie, N. van Dijk, 2011.

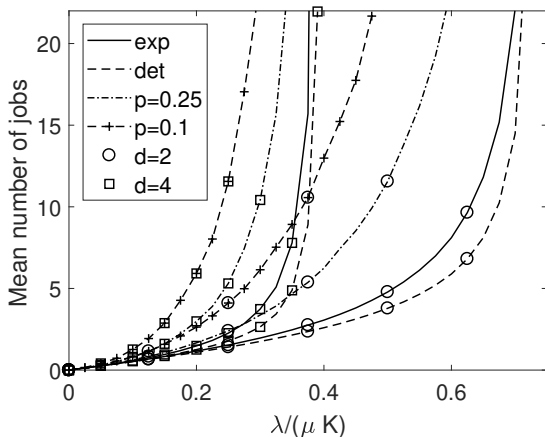
Non-exponential service requirements: FCFS

Mean number of jobs with identical copies and exponential, deterministic and degenerate hyperexponential service requirements.



Non-exponential service requirements: FCFS

Mean number of jobs with identical copies and exponential, deterministic and degenerate hyperexponential service requirements.



Approximations developed by: I Adan, M. Boon, G. Weiss, work in progress

What about heterogeneous systems? ¹⁹

Assume $\mu_1 < \dots < \mu_K$, redundancy- d .

- ▶ With redundancy, the stability condition is

$$\lambda^R = \min_{i=d, \dots, K} \left\{ \mu_i \frac{\binom{K}{d}}{\binom{i-1}{d-1}} \right\}.$$

¹⁹E. Anton, U. Ayesta, M. Jonckheere, I. M. Verloop. Improving the Performance of Heterogeneous Data Centers through Redundancy, to appear in ACM SIGMETRICS 2021

What about heterogeneous systems? ¹⁹

Assume $\mu_1 < \dots < \mu_K$, redundancy- d .

- ▶ With redundancy, the stability condition is

$$\lambda^R = \min_{i=d, \dots, K} \left\{ \mu_i \frac{\binom{K}{d}}{\binom{d-1}{d-1}} \right\}.$$

- ▶ With Bernoulli routing, the stability condition is $\lambda^B = K\mu_1$.

¹⁹E. Anton, U. Ayesta, M. Jonckheere, I. M. Verloop. Improving the Performance of Heterogeneous Data Centers through Redundancy, to appear in ACM SIGMETRICS 2021

What about heterogeneous systems? ¹⁹

Assume $\mu_1 < \dots < \mu_K$, redundancy- d .

- ▶ With redundancy, the stability condition is

$$\lambda^R = \min_{i=d, \dots, K} \left\{ \mu_i \frac{\binom{K}{d}}{\binom{i-1}{d-1}} \right\}.$$

- ▶ With Bernoulli routing, the stability condition is $\lambda^B = K\mu_1$.

Redundancy- d has larger stability region than Bernoulli if $\mu_1 d < \mu_d$.

¹⁹E. Anton, U. Ayesta, M. Jonckheere, I. M. Verloop. Improving the Performance of Heterogeneous Data Centers through Redundancy, to appear in ACM SIGMETRICS 2021

Relation with matching models

Model A:²⁰ :Arriving jobs wait in the queue for a compatible server, arriving servers match the first compatible job and leave the system (even if unmatched)



²⁰I. Adan and R. Righter and G. Weiss, FCFS Parallel Service Systems and Matching Models, Valuetools 2017

²¹I. Adan, A. Busic, J. Mairesse, and G. Weiss. 2015. Reversibility and further properties of FCFS infinite bipartite matching. Math. Oper. Res. 40(2): 500-521 (2015)

Relation with matching models

Model A:²⁰ :Arriving jobs wait in the queue for a compatible server, arriving servers match the first compatible job and leave the system (even if unmatched)



Same distribution as OI queues

²⁰I. Adan and R. Righter and G. Weiss, FCFS Parallel Service Systems and Matching Models, Valuetools 2017

²¹I. Adan, A. Busic, J. Mairesse, and G. Weiss. 2015. Reversibility and further properties of FCFS infinite bipartite matching. Math. Oper. Res. 40(2), 500-521 (2015).

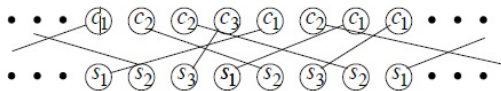
Relation with matching models

Model A:²⁰ :Arriving jobs wait in the queue for a compatible server, arriving servers match the first compatible job and leave the system (even if unmatched)



Same distribution as OI queues

Model B:²¹ (FCFS infinite bipartite matching): arrivals are job server pairs, both of which can queue; arriving jobs (servers) match the first compatible server (job) if any, otherwise join queue



²⁰I. Adan and R. Righter and G. Weiss, FCFS Parallel Service Systems and Matching Models, Valuetools 2017

²¹I. Adan, A. Busic, J. Mairesse, and G. Weiss. 2015. Reversibility and further properties of FCFS infinite bipartite matching. Math. Oper. Res. 40(2): 500-521 (2015)

Conclusions and Perspectives

- ▶ Establish relation between multi-server systems with central queue: redundancy, JSW, etc

Conclusions and Perspectives

- ▶ Establish relation between multi-server systems with central queue: redundancy, JSW, etc
- ▶ Token based and Gen-OI provide unifying framework to cover OI and “multi-job multi-server” models

Conclusions and Perspectives

- ▶ Establish relation between multi-server systems with central queue: redundancy, JSW, etc
- ▶ Token based and Gen-OI provide unifying framework to cover OI and “multi-job multi-server” models
- ▶ Product form distribution does not directly imply computability
⇒ computation of $\pi(0)$.

Conclusions and Perspectives

- ▶ Establish relation between multi-server systems with central queue: redundancy, JSW, etc
- ▶ Token based and Gen-OI provide unifying framework to cover OI and “multi-job multi-server” models
- ▶ Product form distribution does not directly imply computability
⇒ computation of $\pi(0)$.
- ▶ Not final word on product form distributions yet...
⇒ Pass & Swap queues (Comte& Dorsman 2020) have product form, yet not included in Token or Gen-OI

Conclusions and Perspectives

- ▶ Most of the analysis on redundancy assumes FCFS
 - ▶ Need to consider other classical disciplines
 - ▶ Develop new redundancy-aware scheduling disciplines

Conclusions and Perspectives

- ▶ Most of the analysis on redundancy assumes FCFS
 - ▶ Need to consider other classical disciplines
 - ▶ Develop new redundancy-aware scheduling disciplines
- ▶ Much less known for non exponential service times

Conclusions and Perspectives

- ▶ Most of the analysis on redundancy assumes FCFS
 - ▶ Need to consider other classical disciplines
 - ▶ Develop new redundancy-aware scheduling disciplines
- ▶ Much less known for non exponential service times
- ▶ Exploit steady-state distribution to characterize heavy-traffic, mean-field results, asymptotic optimality etc.
 - ⇒ E. Cardinaels, S. Borst, J.H. van Leeuwen, Redundancy Scheduling with Locally Stable Compatibility Graphs, arxiv 2020

Conclusions and Perspectives

- ▶ Most of the analysis on redundancy assumes FCFS
 - ▶ Need to consider other classical disciplines
 - ▶ Develop new redundancy-aware scheduling disciplines
- ▶ Much less known for non exponential service times
- ▶ Exploit steady-state distribution to characterize heavy-traffic, mean-field results, asymptotic optimality etc.
 - ⇒ E. Cardinaels, S. Borst, J.H. van Leeuwen, Redundancy Scheduling with Locally Stable Compatibility Graphs, arxiv 2020
- ▶ Relation between JSQ(d) and Redundancy
 - ⇒ Coming up SNAPP talk by S. Borst

(Partial) Bibliography

Redundancy:

K. Gardner, S. Zbarsky, S. Doroudi, M. Harchol-Balter, E. Hyttiä, and A. Scheller-Wolf. 2016. Queueing with redundant requests: exact analysis. *Queueing Systems* 83, 3-4 (2016), 227-259.

K. Gardner, M. Harchol-Balter, A. Scheller-Wolf, M. Velednitsky, S. Zbarsky, *Redundancy-d: The power of d choices for redundancy*, *Operations Research*, 2017

Multi-type job and servers:

Visschers, Adan, Weiss, *A product form solution to a system with multi-type jobs and multi-type servers*, *Queueing Systems*, 2012.

Adan, Weiss, *A skill based parallel service system under FCFS-ALIS: steady state, overloads, and abandonments* *Stochastic Systems*, *INFORMS*, 2014, 4, 250-299

OI Queues:

A. Krzesinski, *Order independent queues*, in “*Queueing Networks: a fundamental approach*”, eds: R. Boucherie, N. van Dijk, 2011.

(Partial) Bibliography

Redundancy and product form

Ayesta, Bodas, Verloop, On a unifying product form framework for redundancy models, Performance Evaluation, 2018

T. Bonald, C. Comte, F Mathieu, Performance of Balanced Fairness in Resource Pools: A Recursive Approach, Sigmetrics 2017

K. Gardner and R. Righter, Product (Re)forms, Tutorial APS 2019, available at: <https://kgardner.people.amherst.edu/>

Mean-Field and Heavy-Traffic:

M. Bramson, Yi Lu, B. Prabhakar, Randomized load balancing with general service time distributions. SIGMETRICS 2010: 275-286

T. Hellemsans, B. Van Houdt, On the Power-of-d-choices with Least Loaded Server Selection, SIGMETRICS 2018

E. Cardinaels, S. Borst, J.H. van Leeuwen, Redundancy Scheduling with Locally Stable Compatibility Graphs, arxiv 2020

Generalizations:

Ayesta, Bodas, Dorsman, Verloop, A token-based central queue with order-independent service rates , to appear in OR

K. Gardner, R. Righter, Product Forms for FCFS Queueing Models with Arbitrary Server-Job Compatibilities: An Overview, to appear in QUESTA

(Partial) Bibliography

Stability:

E. Anton, U. Ayesta, M. Jonckheere, I. M. Verloop. On the stability of redundancy models, to appear in OR

Y. Raaijmakers, S. Borst, O. Boxma, Delta probing policies for redundancy. Perform. Eval (2018)

E. Anton, U. Ayesta, M. Jonckheere, I. M. Verloop. Improving the Performance of Heterogeneous Data Centers through Redundancy, to appear in ACM SIGMETRICS 2021

Scheduling and efficiency in redundancy:

K. Gardner, M. Harchol-Balter, E. Hyttiä, R. Righter: Scheduling for efficiency and fairness in systems with redundancy. Perform. Evaluation 116: 1-25 (2017)

Y. Raaijmakers, S. Borst, O. Boxma, Delta probing policies for redundancy. Perform. Eval (2018)

Matching:

I. Adan and R. Righter and G. Weiss, FCFS Parallel Service Systems and Matching Models, Valuetools 2017

I. Adan, A. Busic, J. Mairesse, and G. Weiss. 2015. Reversibility and further properties of FCFS infinite bipartite matching. Math. Oper. Res. 43(2): 598-621 (2018)

J. Mairesse, P. Moyal: Stability of the stochastic matching model. J. Appl. Probab. 2016

Load balancing, redundancy, and multi type job and server systems

Urtzi Ayesta

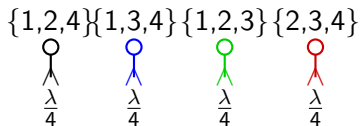
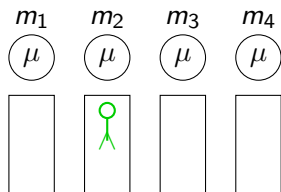
CNRS & Ikerbasque-Univ. Basque Country

Questions/Remarks welcome at : urtzi.ayesta@irit.fr

IFIP Performance 2020, 2/11/2020

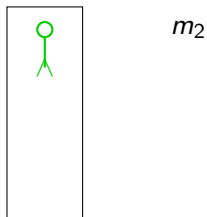
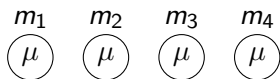
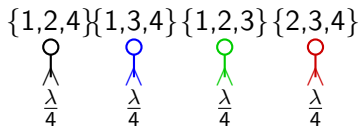
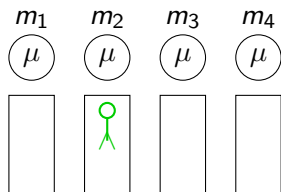
redundancy- d with c.o.s.

$$N = 4 \text{ and } d = 3$$



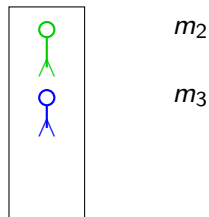
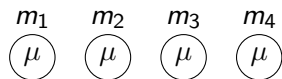
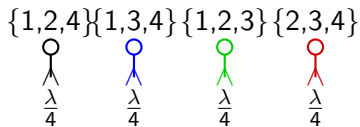
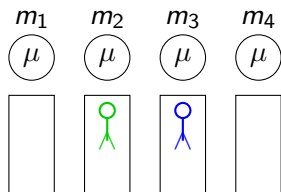
redundancy- d with c.o.s.

$N = 4$ and $d = 3$



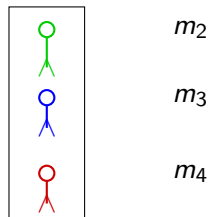
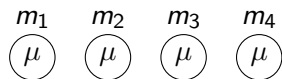
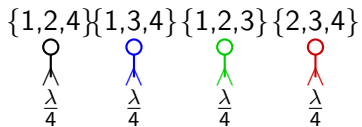
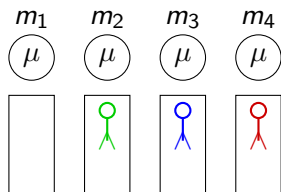
redundancy- d with c.o.s.

$N = 4$ and $d = 3$



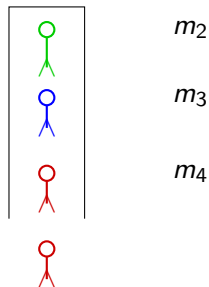
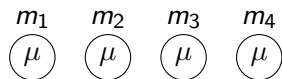
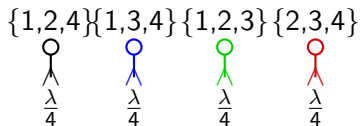
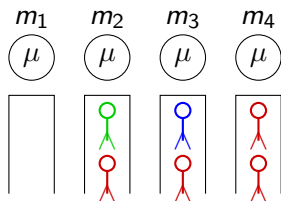
redundancy- d with c.o.s.

$N = 4$ and $d = 3$



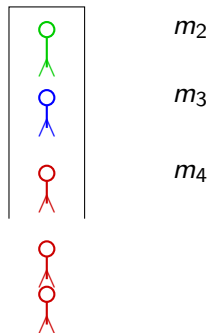
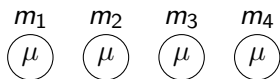
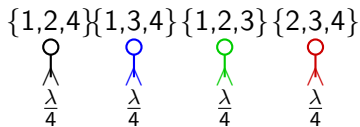
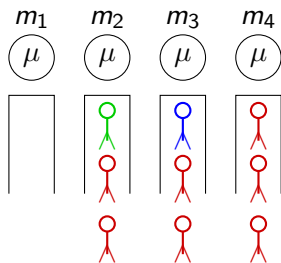
redundancy- d with c.o.s.

$N = 4$ and $d = 3$



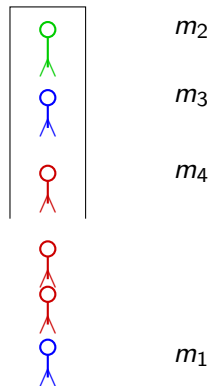
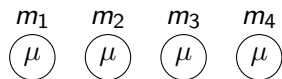
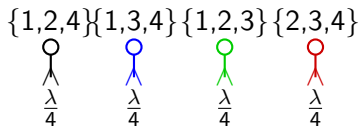
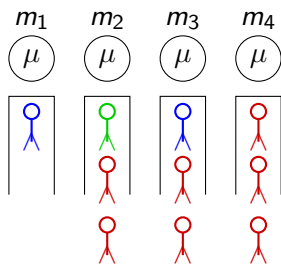
redundancy- d with c.o.s.

$N = 4$ and $d = 3$



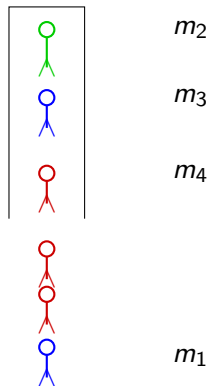
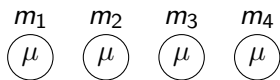
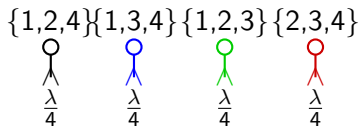
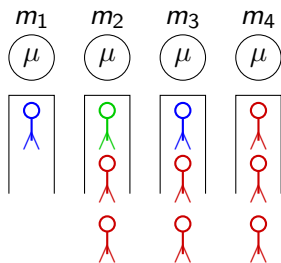
redundancy- d with c.o.s.

$N = 4$ and $d = 3$



redundancy- d with c.o.s.

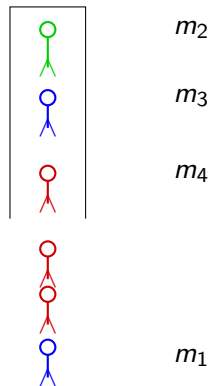
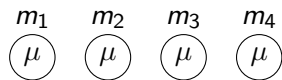
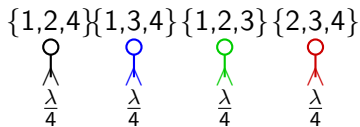
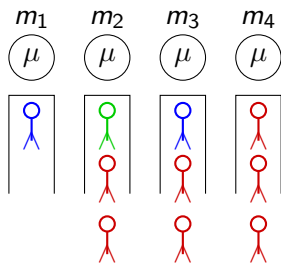
$N = 4$ and $d = 3$



$(m_2, m_3, m_4, 2, m_1)$

redundancy- d with c.o.s.

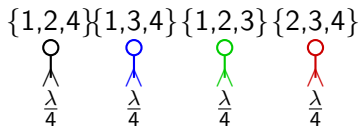
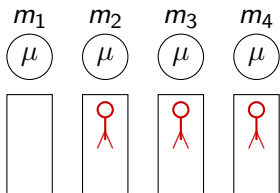
$N = 4$ and $d = 3$



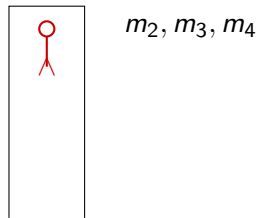
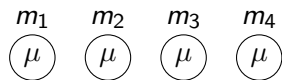
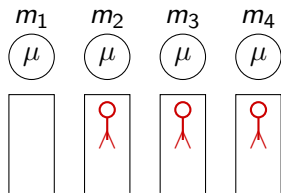
$(m_2, m_3, m_4, 2, m_1)$

$(M_1, \dots, M_d, n_d, M_{d+1}, \dots, n_i, M_i)$

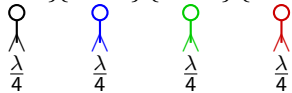
redundancy- d with c.o.c.



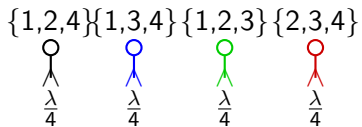
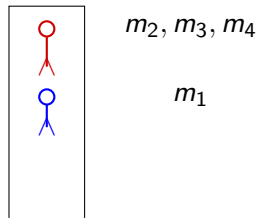
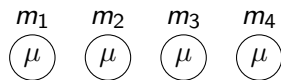
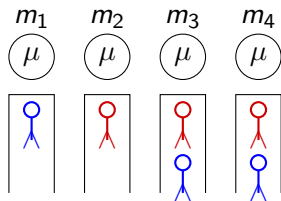
redundancy- d with c.o.c.



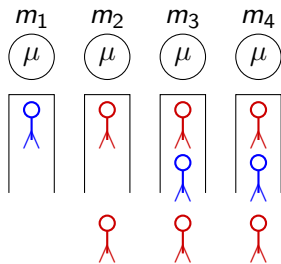
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



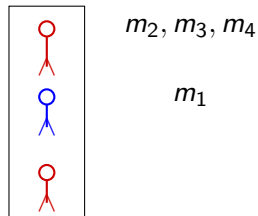
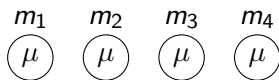
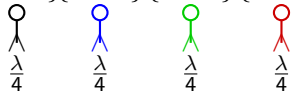
redundancy- d with c.o.c.



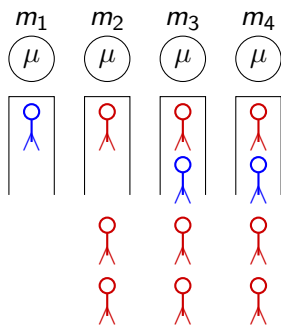
redundancy- d with c.o.c.



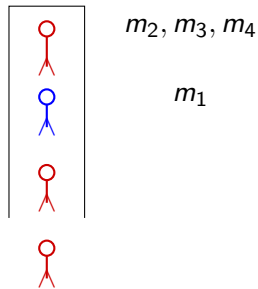
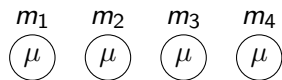
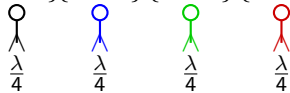
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



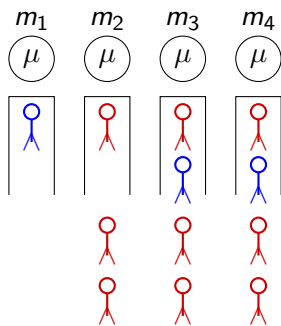
redundancy- d with c.o.c.



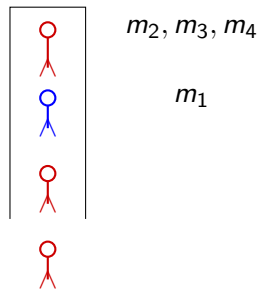
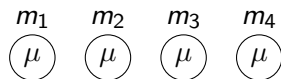
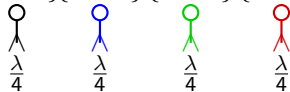
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



redundancy- d with c.o.c.

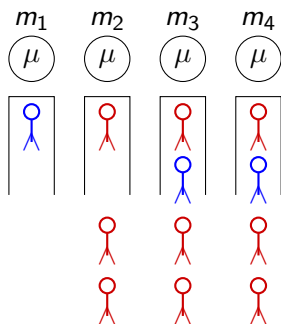


$\{1,2,4\} \{1,3,4\} \{1,2,3\} \{2,3,4\}$

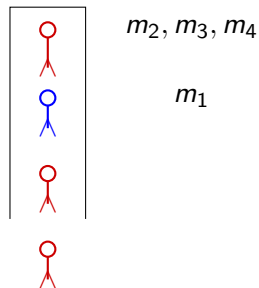
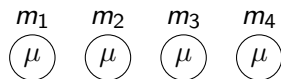
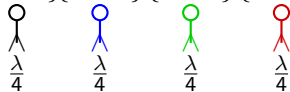


$(F_{red}, F_{blue}, 2)$

redundancy- d with c.o.c.



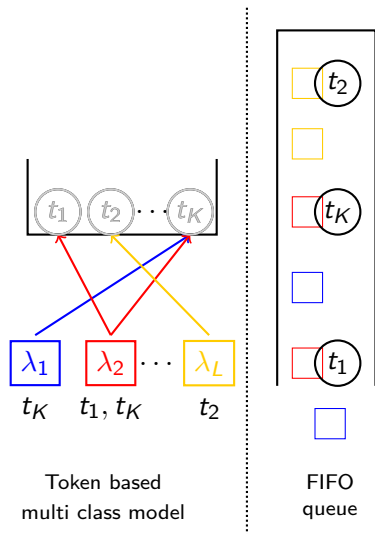
$\{1,2,4\}$ $\{1,3,4\}$ $\{1,2,3\}$ $\{2,3,4\}$



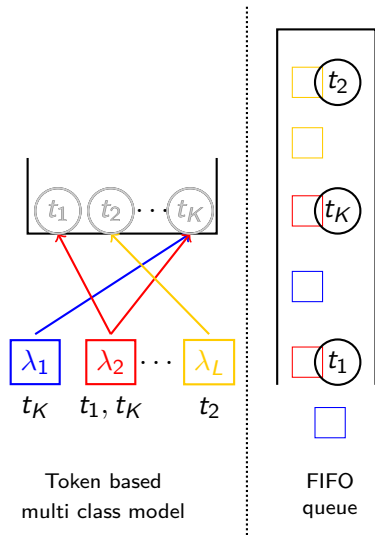
$(F_{red}, F_{blue}, 2)$

$F_{c_1}, n_1, F_{c_2}, n_2, \dots, F_{c_i}, n_i$

Classes that wait

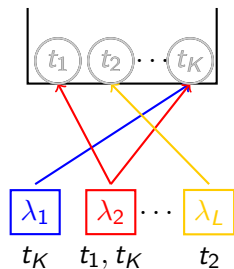


Classes that wait

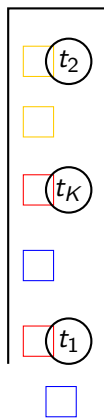


- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$

Classes that wait



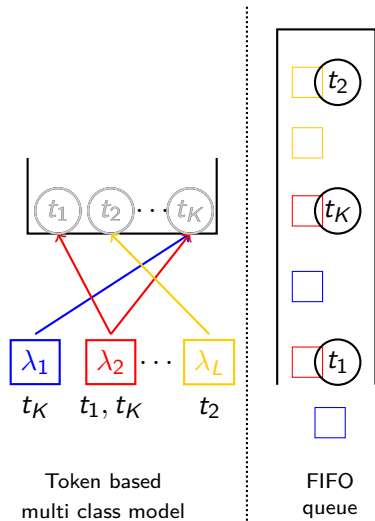
Token based
multi class model



FIFO
queue

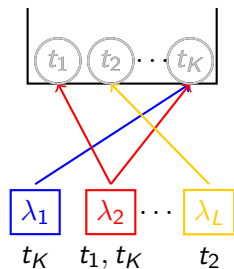
- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- What are the job classes that constitute the n_j jobs ?

Classes that wait

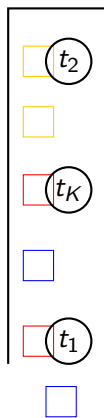


- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- What are the job classes that constitute the n_j jobs ?
- $\mathcal{U}(\{T_1, T_2, \dots, T_j\})$ denotes the set of job classes who have to wait because tokens $\{T_1, T_2, \dots, T_j\}$ are busy

Classes that wait



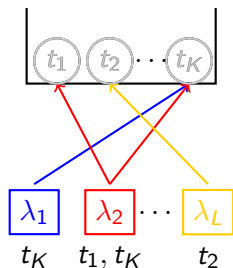
Token based
multi class model



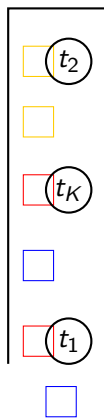
FIFO
queue

- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- What are the job classes that constitute the n_j jobs ?
- $\mathcal{U}(\{T_1, T_2, \dots, T_j\})$ denotes the set of job classes who have to wait because tokens $\{T_1, T_2, \dots, T_j\}$ are busy
- $\mathcal{U}(\{t_2\}) = \{L\}, \mathcal{U}(\{t_2, t_K\}) = \{1, L\}$ and $\mathcal{U}(\{t_2, t_K, t_1\}) = \{1, 2, L\}$

Classes that wait



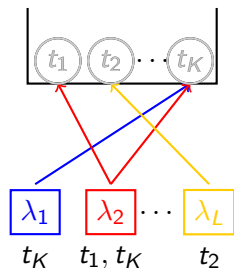
Token based
multi class model



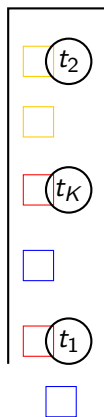
FIFO
queue

- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- What are the job classes that constitute the n_j jobs ?
- $\mathcal{U}(\{T_1, T_2, \dots, T_j\})$ denotes the set of job classes who have to wait because tokens $\{T_1, T_2, \dots, T_j\}$ are busy
- $\mathcal{U}(\{t_2\}) = \{L\}$, $\mathcal{U}(\{t_2, t_K\}) = \{1, L\}$ and $\mathcal{U}(\{t_2, t_K, t_1\}) = \{1, 2, L\}$
- $\lambda_{\mathcal{U}(\{t_2, t_K, t_1\})} = \lambda_1 + \lambda_2 + \lambda_L$

Classes that wait



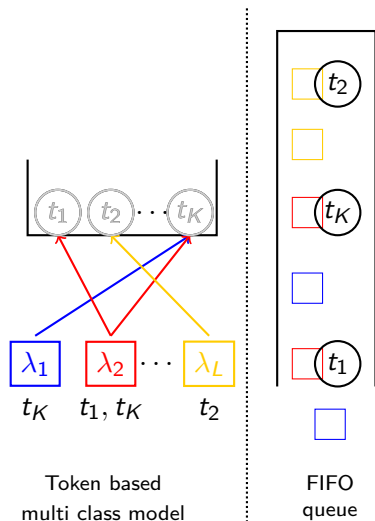
Token based
multi class model



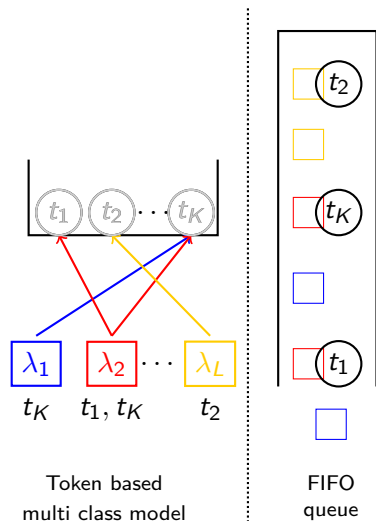
FIFO
queue

- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- What are the job classes that constitute the n_j jobs ?
- $\mathcal{U}(\{T_1, T_2, \dots, T_j\})$ denotes the set of job classes who have to wait because tokens $\{T_1, T_2, \dots, T_j\}$ are busy
- $\mathcal{U}(\{t_2\}) = \{L\}$, $\mathcal{U}(\{t_2, t_K\}) = \{1, L\}$ and $\mathcal{U}(\{t_2, t_K, t_1\}) = \{1, 2, L\}$
- $\lambda_{\mathcal{U}(\{t_2, t_K, t_1\})} = \lambda_1 + \lambda_2 + \lambda_L$
- $\lambda_{\mathcal{U}(\{T_1, T_2, \dots, T_j\})}$

Activating a token

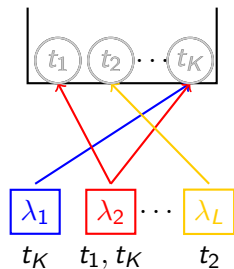


Activating a token

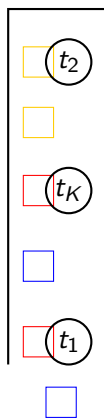


- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$

Activating a token



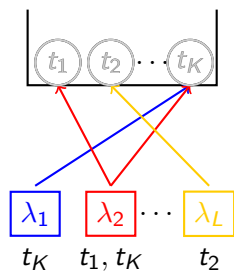
Token based
multi class model



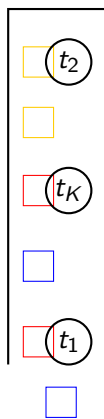
FIFO
queue

- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- Given that tokens $\{T_1, T_2, \dots, T_i\}$ are busy, at what rate does a free token T_{i+1} become busy ?

Activating a token



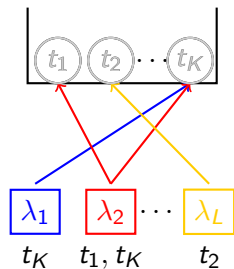
Token based
multi class model



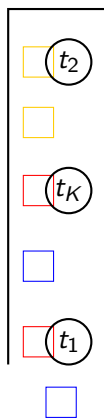
FIFO
queue

- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- Given that tokens $\{T_1, T_2, \dots, T_i\}$ are busy, at what rate does a free token T_{i+1} become busy ?
- Denoted by $\lambda_{T_{i+1}}(\{T_1, T_2, \dots, T_i\})$

Activating a token



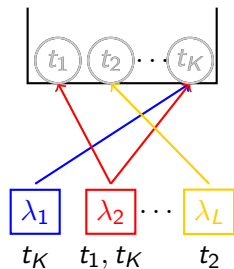
Token based
multi class model



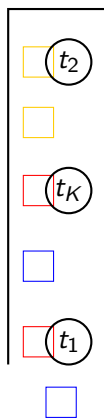
FIFO
queue

- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- Given that tokens $\{T_1, T_2, \dots, T_i\}$ are busy, at what rate does a free token T_{i+1} become busy ?
- Denoted by $\lambda_{T_{i+1}}(\{T_1, T_2, \dots, T_i\})$
- $\lambda_{t_2}(\{\phi\}) = \lambda_L$

Activating a token



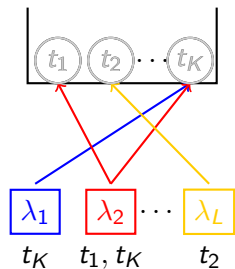
Token based
multi class model



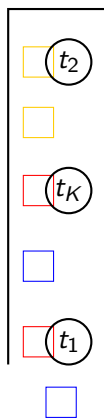
FIFO
queue

- Consider a generic state $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- Given that tokens $\{T_1, T_2, \dots, T_i\}$ are busy, at what rate does a free token T_{i+1} become busy ?
- Denoted by $\lambda_{T_{i+1}}(\{T_1, T_2, \dots, T_i\})$
- $\lambda_{t_2}(\{\phi\}) = \lambda_L$
- $\lambda_{t_K}(\{t_2\}) = \lambda_1 + p\lambda_2$

Activating a token



Token based
multi class model



FIFO
queue

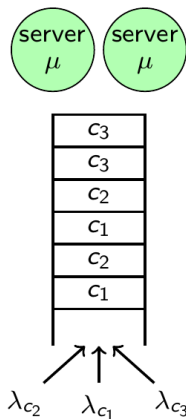
- Consider a generic state
 $x = (T_1, n_1, \dots, T_j, n_j, \dots, T_i, n_i)$
- Given that tokens $\{T_1, T_2, \dots, T_i\}$ are busy, at what rate does a free token T_{i+1} become busy ?
- Denoted by $\lambda_{T_{i+1}}(\{T_1, T_2, \dots, T_i\})$
- $\lambda_{t_2}(\{\phi\}) = \lambda_L$
- $\lambda_{t_K}(\{t_2\}) = \lambda_1 + \rho\lambda_2$
- $\lambda_{t_1}(\{t_2, t_K\}) = \lambda_2$

Applications of OI

- ▶ Multi-class M/M/1 queue

Applications of OI

- ▶ Multi-class M/M/1 queue
- ▶ MSCCC - Multi-server with concurrent customers

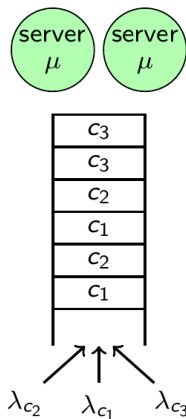


2 identical servers, 3 classes,
servers compatible with all classes

Concurrent control: Both servers
cannot serve the same class simul-
taneously

Applications of OI

- ▶ Multi-class M/M/1 queue
- ▶ MSCCC - Multi-server with concurrent customers



2 identical servers, 3 classes,
servers compatible with all classes

Concurrent control: Both servers
cannot serve the same class simul-
taneously

- ▶ Processor-Sharing systems: $\mu(c_1, \dots, c_n)$ might depend on a scalar function $\phi(n)$