



Waiting in an efficient way

Scalable algorithms to dispatch jobs efficiently in large-scale many-server queueing networks.

Mark van der Boor

March, 2019

What are we talking about?

Queueing Theory



A busy Sunday in a supermarket. Foto: Harold Röring

Waiting in an efficient way

What are we talking about?

Queueing Theory
Decision Theory



A busy Sunday in a supermarket. Foto: Harold Röring

Waiting in an efficient way

What are we talking about?

Queueing Theory
Decision Theory
Stochastics



A busy Sunday in a supermarket. Foto: Harold Röring

Waiting in an efficient way

What are we talking about?

Queueing Theory

Decision Theory

Stochastics

“Load Balancing”



A busy Sunday in a supermarket. Foto: Harold Röring



Waiting in an efficient way



Waiting in an efficient way



Waiting in an efficient way



Waiting in an efficient way



Waiting in an efficient way



Waiting in an efficient way

Load Balancing

Jobs arrive at a central dispatcher
and need to be send to one of many servers .

Load Balancing

People arrive at the queues in a shop and need to be send to one of many cashiers .



Load Balancing

Texts arrive at a central server and need to be sent to one of many satellites.



Load Balancing

Cars arrive at a toll road and need to be send to one of many booths .



Load Balancing

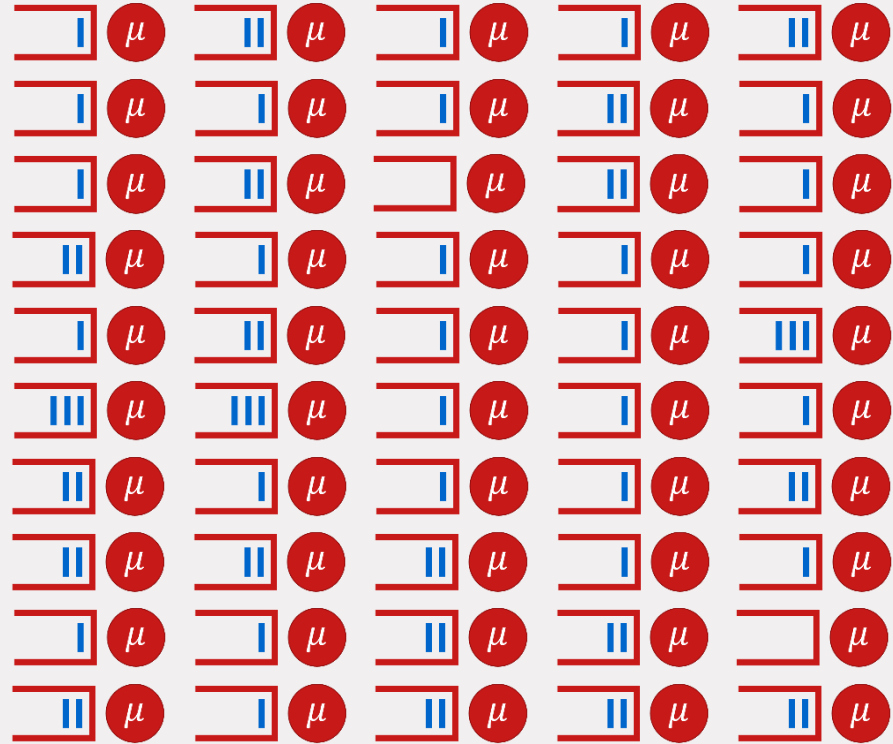
Data arrives at a dispatcher and need to be send to one of many servers .



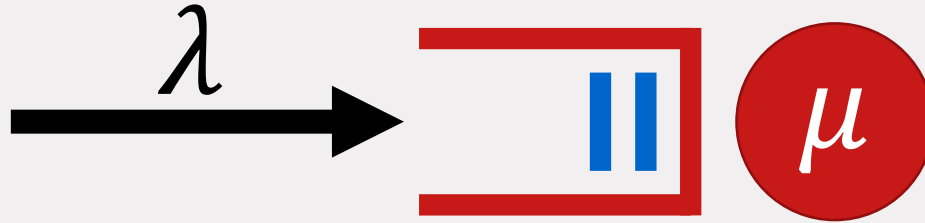
Waiting in an efficient way

Mathematical model

- Model every server as a queue
- For every incoming job, route it to one of the queues.
- What is the smartest you can do?

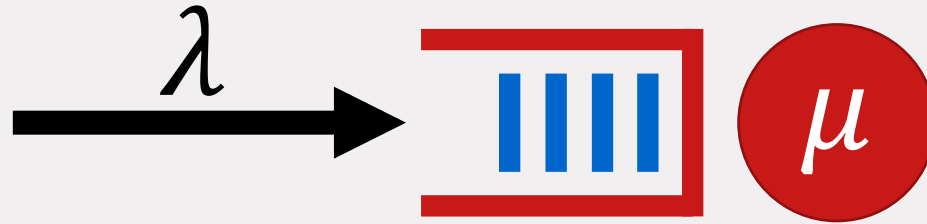


Mathematics time!*



*Actually, modelling a problem is also (a very important) part of mathematics!

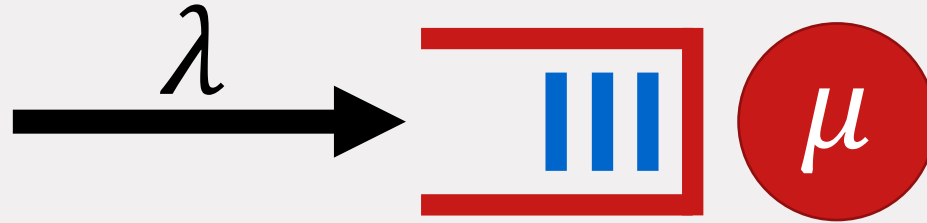
Mathematics time!



Number of jobs in the system:



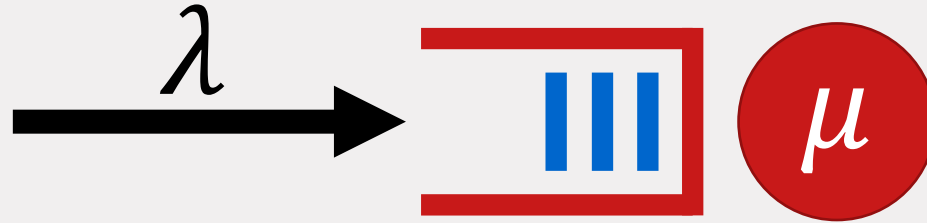
Mathematics time!



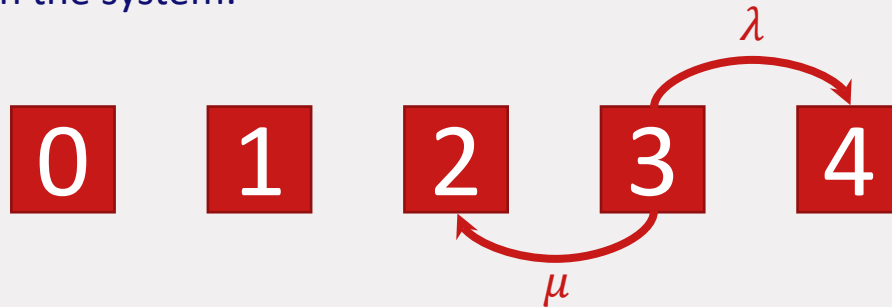
Number of jobs in the system:



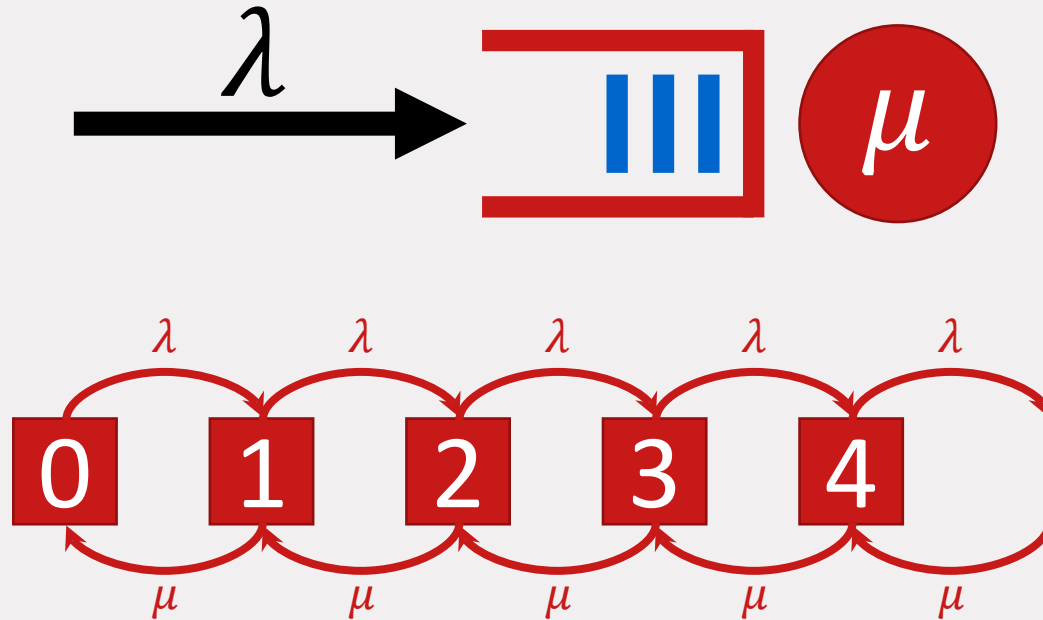
Mathematics time!



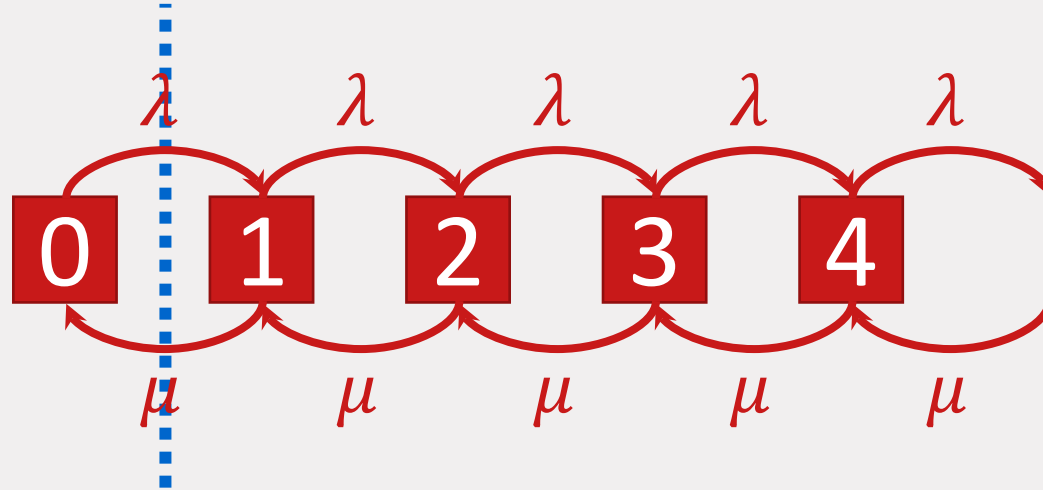
Number of jobs in the system:



Mathematics time!



Mathematics time!



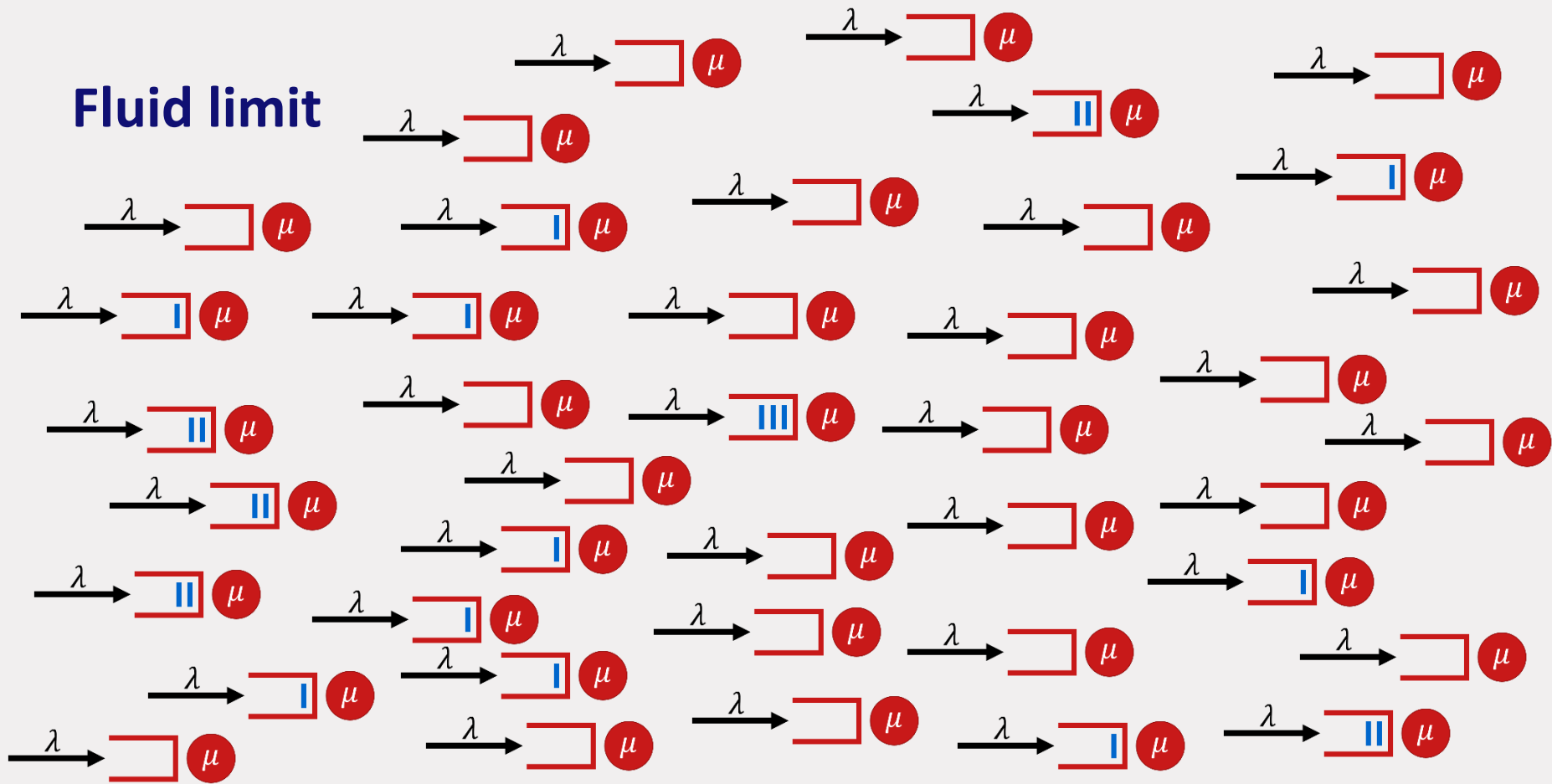
Waiting in an efficient way

Exercise – M | M | c queue

In the previous example, only the first job in line gets service. What would happen if c jobs can get served simultaneously?

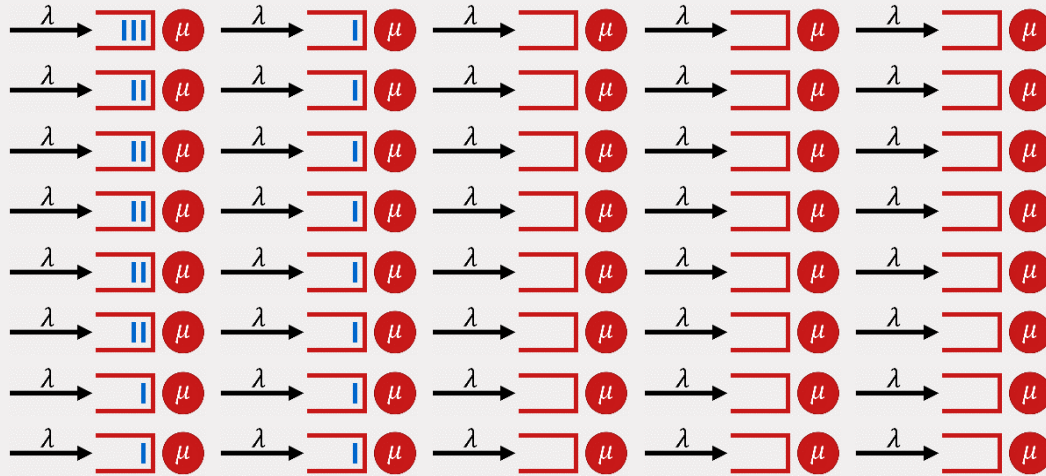
- a) Draw the transition diagram.
- b) Can you find the equilibrium distribution?

Fluid limit



Waiting in an efficient way

Fluid limit



Let's look at the fraction of systems that have i jobs at time 0. We'll call this $f_i(0)$.

In this case:

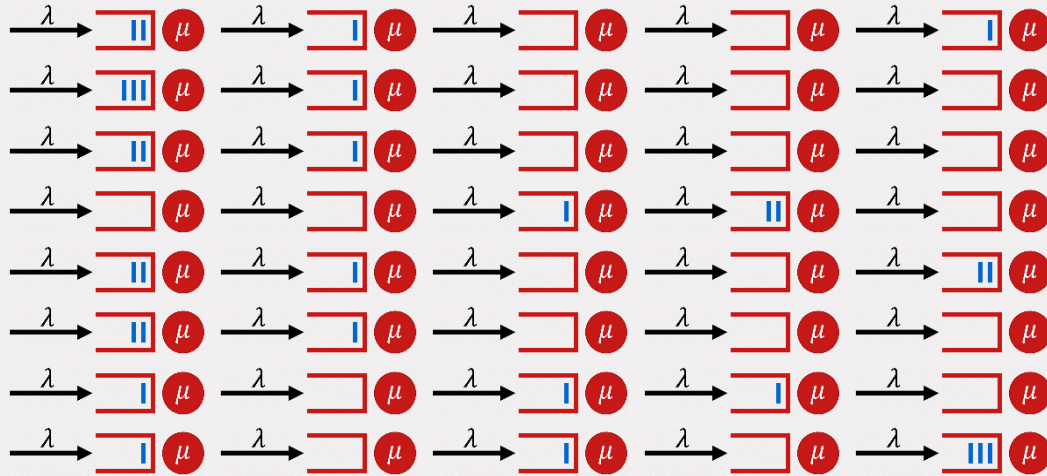
$$f_0(0) = 0.60$$

$$f_1(0) = 0.25$$

$$f_2(0) = 0.125$$

$$f_3(0) = 0.025$$

Fluid limit



At a later time t , these fractions will have changed.

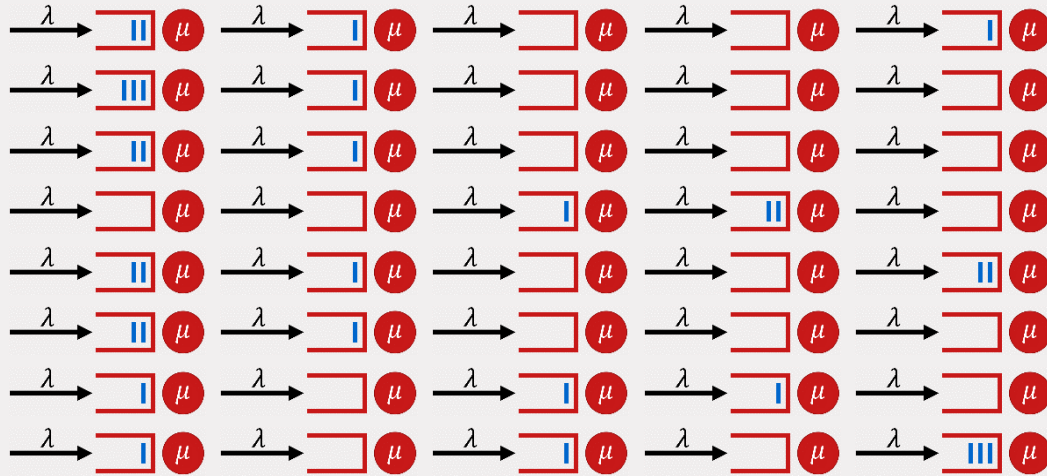
$$f_0(t) = 0.50$$

$$f_1(t) = 0.30$$

$$f_2(t) = 0.15$$

$$f_3(t) = 0.05$$

Fluid limit

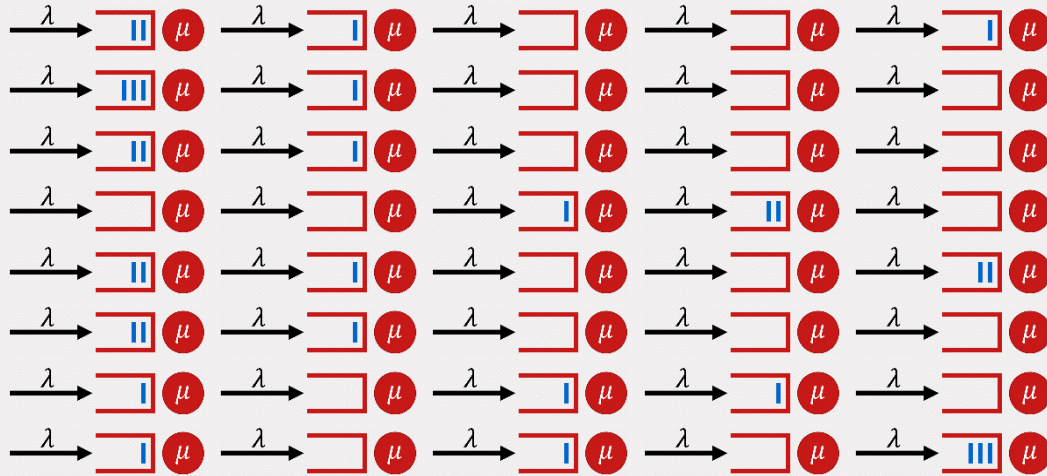


Let's look at how quickly these fractions change

(we assume we have infinitely many servers)

$$\frac{df_0(t)}{dt} = -\lambda \cdot f_0(t)$$

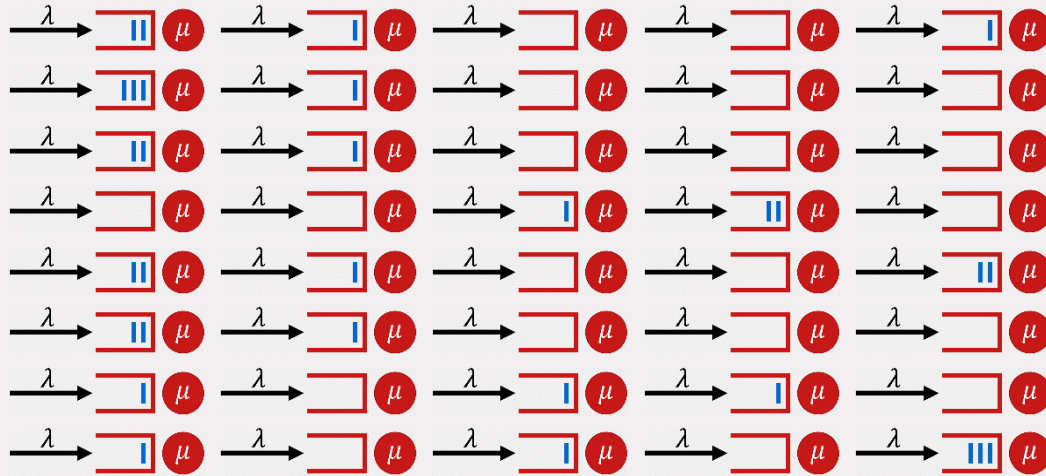
Fluid limit



Let's look at how quickly these fractions change

$$\frac{df_0(t)}{dt} = -\lambda \cdot f_0(t) + \mu \cdot f_1(t)$$

Fluid limit



Let's look at how quickly these fractions change

$$\frac{df_0(t)}{dt} = -\lambda \cdot f_0(t) + \mu \cdot f_1(t)$$

Differential Equations

$$\frac{df_0(t)}{dt} = -\lambda \cdot f_0(t) + \mu \cdot f_1(t)$$

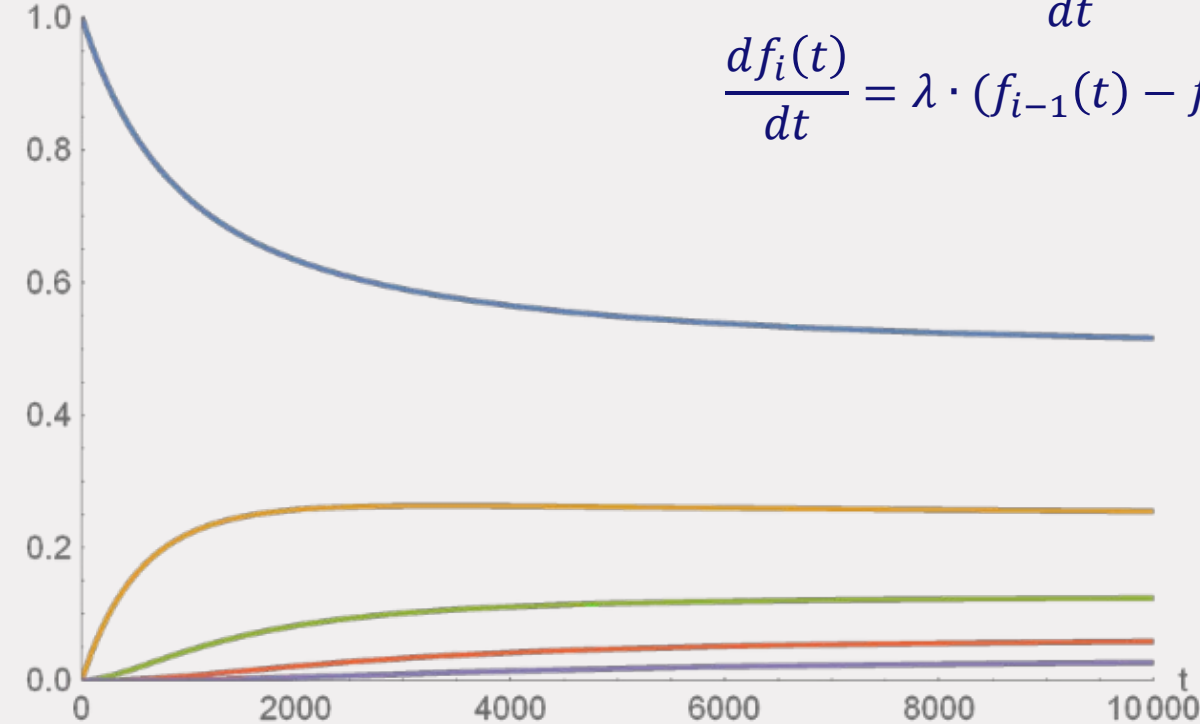
$$\frac{df_i(t)}{dt} = \lambda \cdot (f_{i-1}(t) - f_i(t)) + \mu \cdot (f_{i+1}(t) - f_i(t)), i \geq 1$$

Differential Equations

$$\frac{df_0(t)}{dt} = -\lambda \cdot f_0(t) + \mu \cdot f_1(t)$$

$$\frac{df_i(t)}{dt} = \lambda \cdot (f_{i-1}(t) - f_i(t)) + \mu \cdot (f_{i+1}(t) - f_i(t)), i \geq 1$$

Fraction

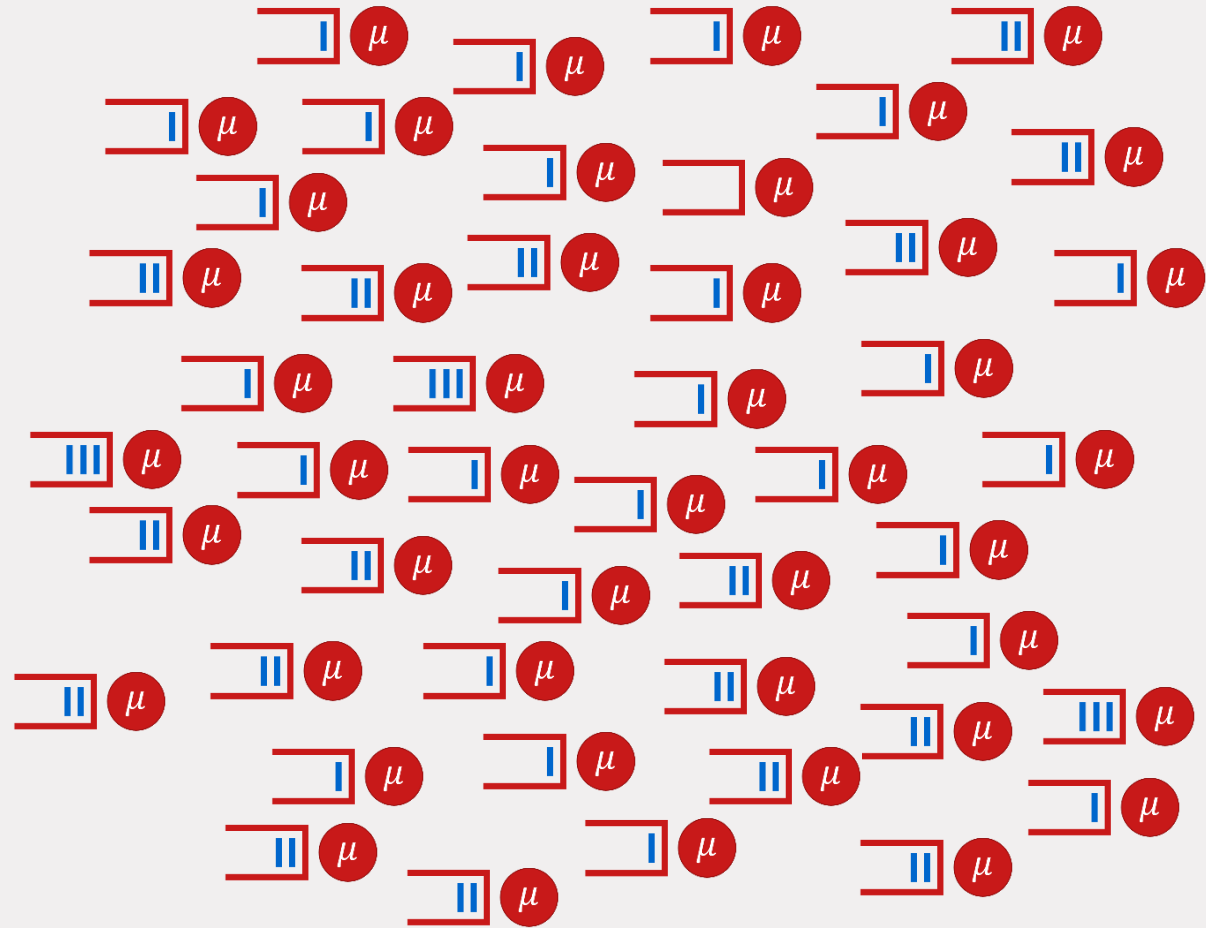


- $f_0(t)$
- $f_1(t)$
- $f_2(t)$
- $f_3(t)$
- $f_4(t)$

Waiting in an efficient way

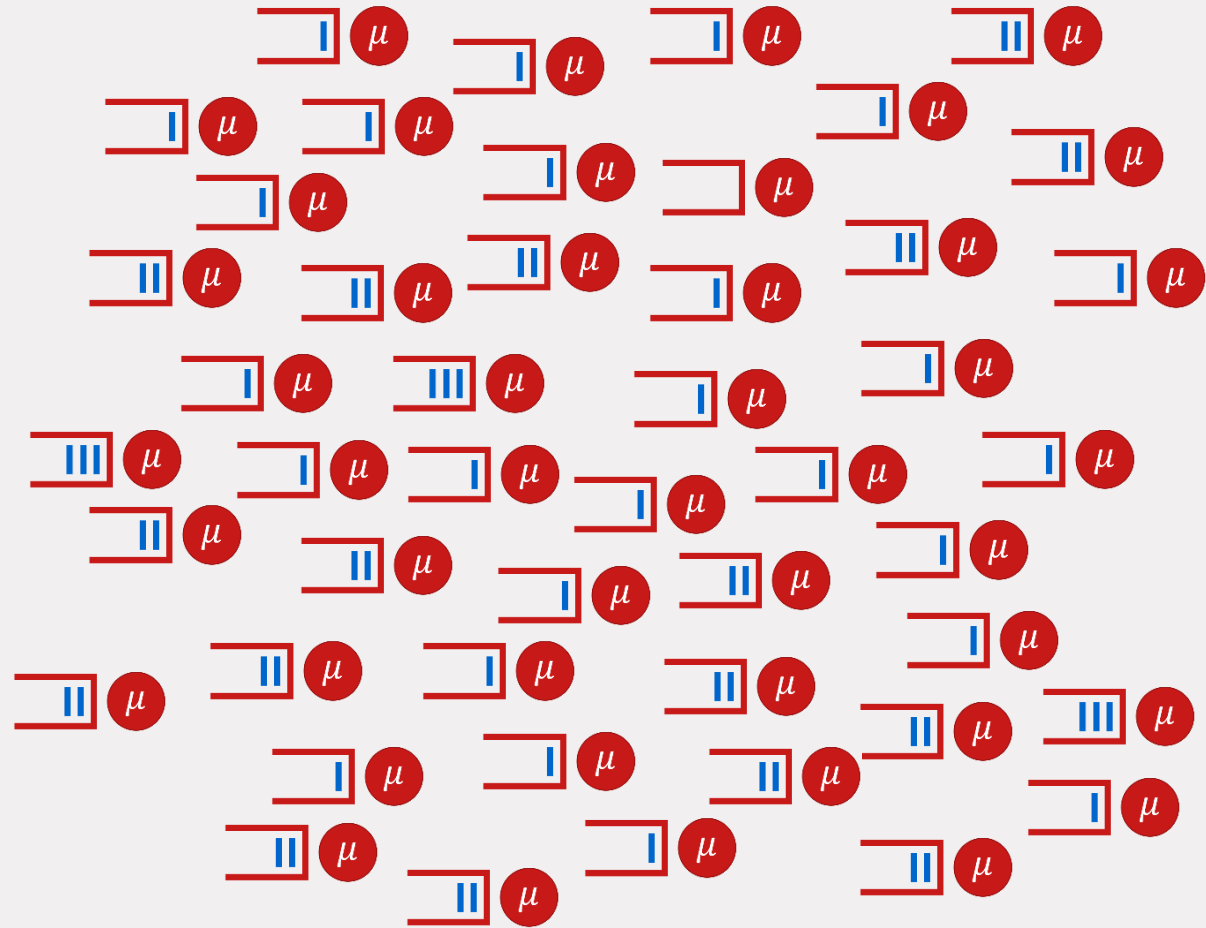
Load balancing

- Join the Shortest Queue
- Random
- Power-of-2



Power-of-2

- Choose two random servers
- Pick the one with the shortest queue



Power-of-2

Define $g_i(t)$ to be the fraction of servers with queue length $\geq i$.

Power-of-2

Define $g_i(t)$ to be the fraction of servers with queue length $\geq i$.

$$\frac{dg_i(t)}{dt} = g_{i+1}(t) - g_i(t) + \lambda(g_{i-1}(t)^2 - g_i(t)^2)$$

Power-of-2

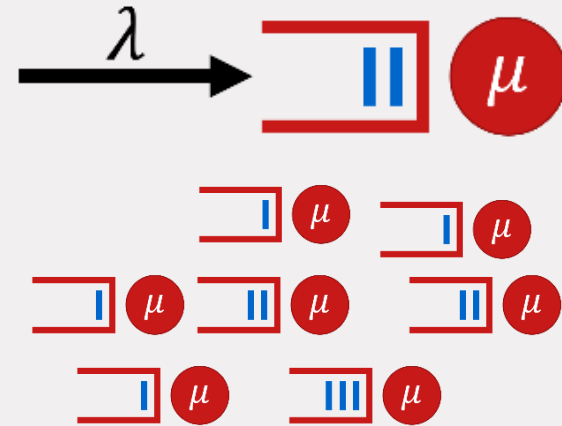
Define $g_i(t)$ to be the fraction of servers with queue length $\geq i$.

$$0 = g_{i+1}^* - g_i^* + \lambda(g_{i-1}^{*2} - g_i^{*2})$$

$$\Rightarrow g_i^* \sim \lambda^{2^i - 1}$$

Simulations





Than-Queue!