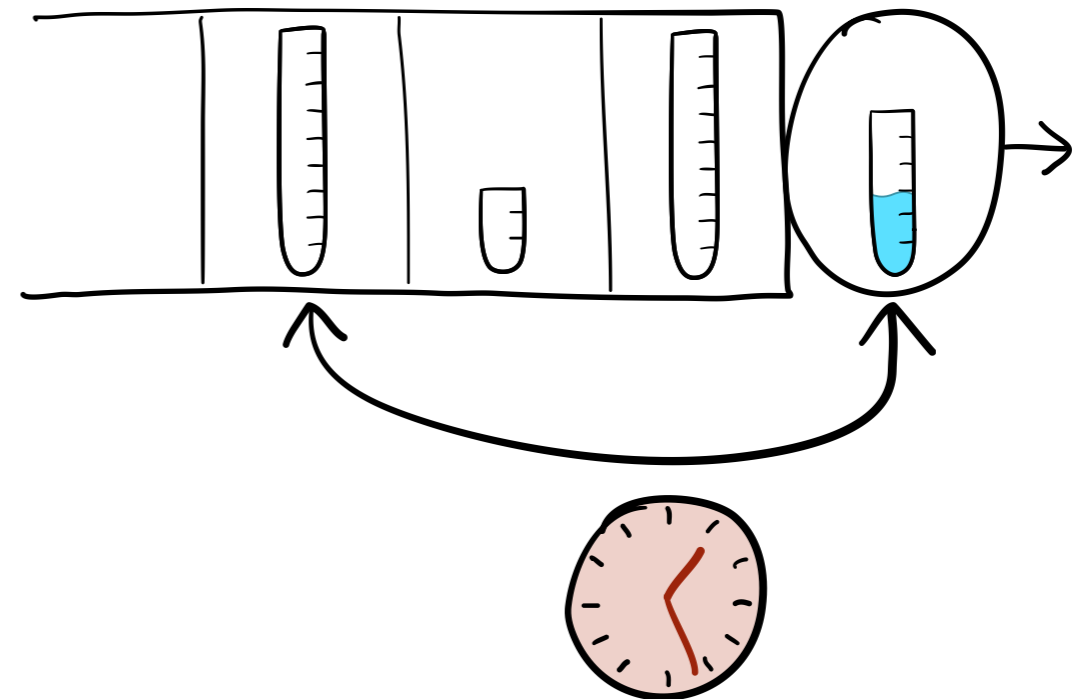


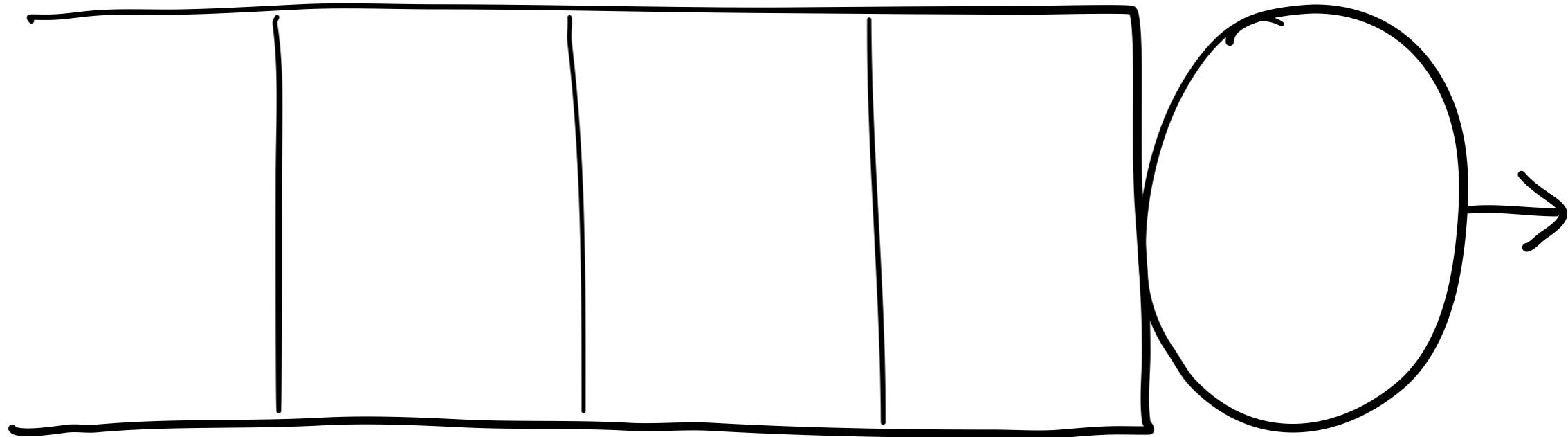
# *Queueing with* **Preemption Costs**



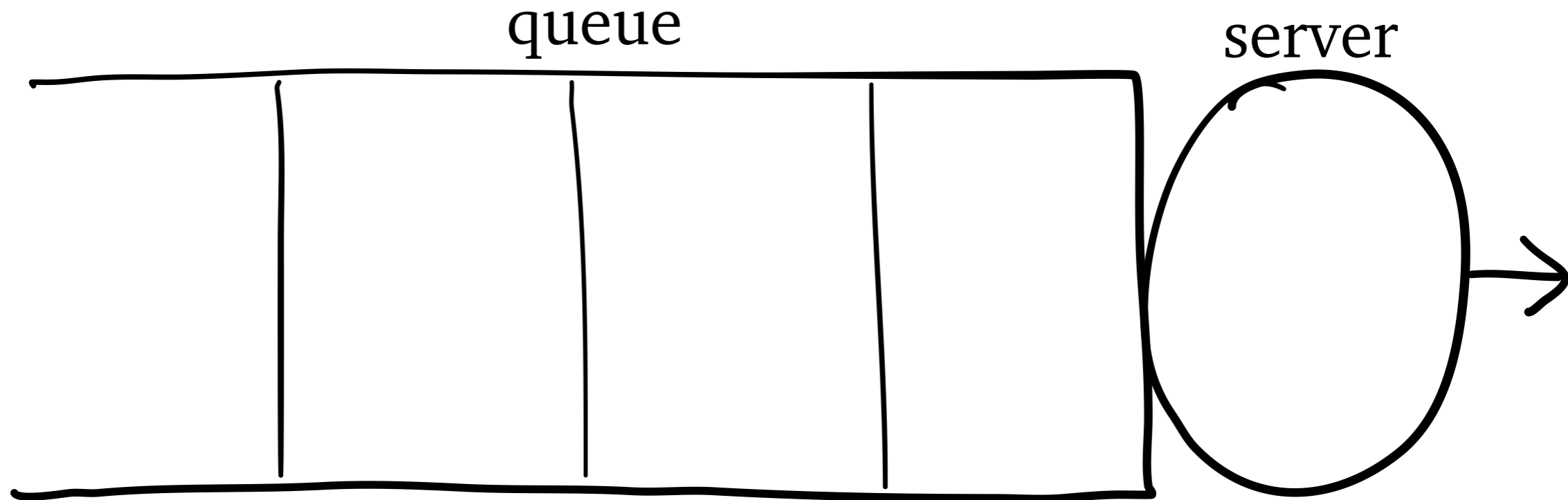
Ziv Scully

Carnegie Mellon University

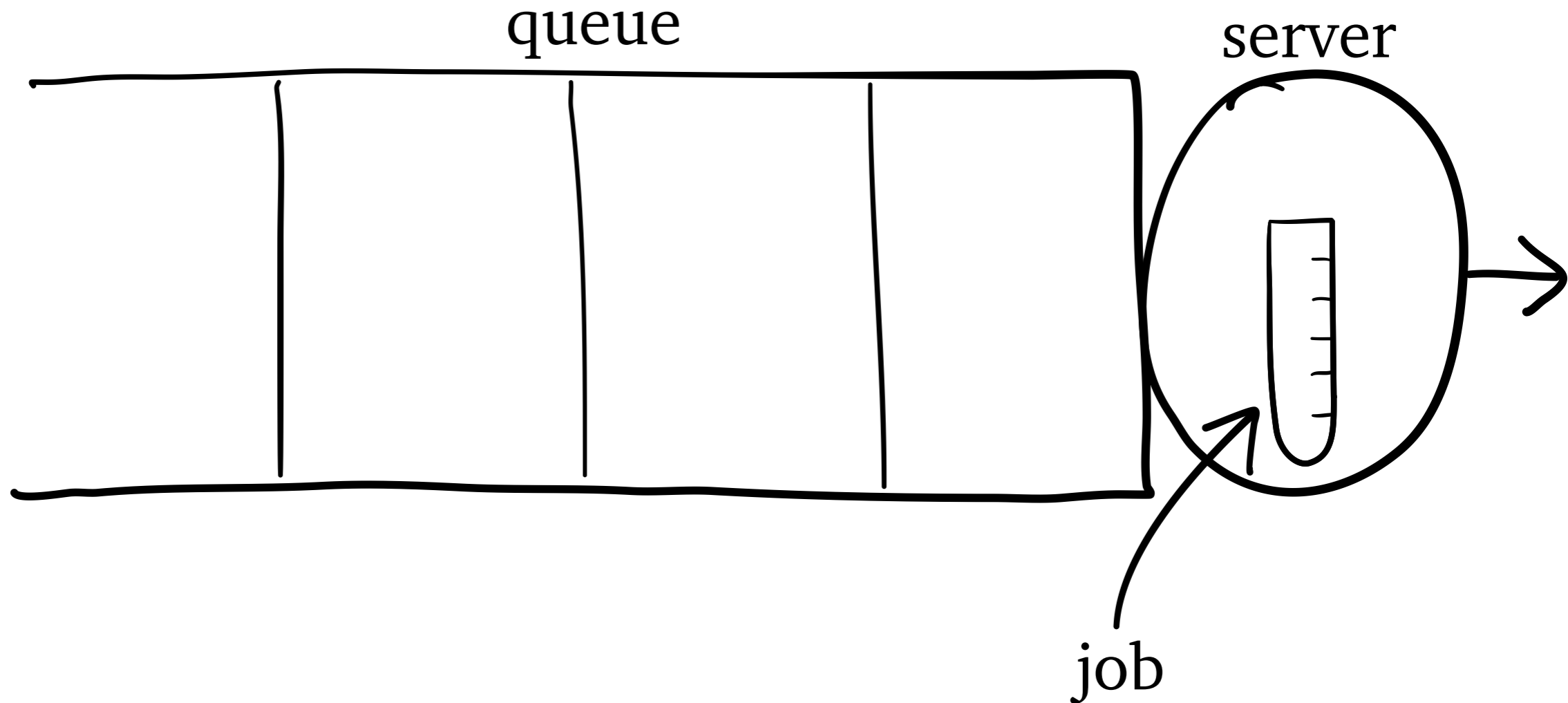
# M/G/1 Queue



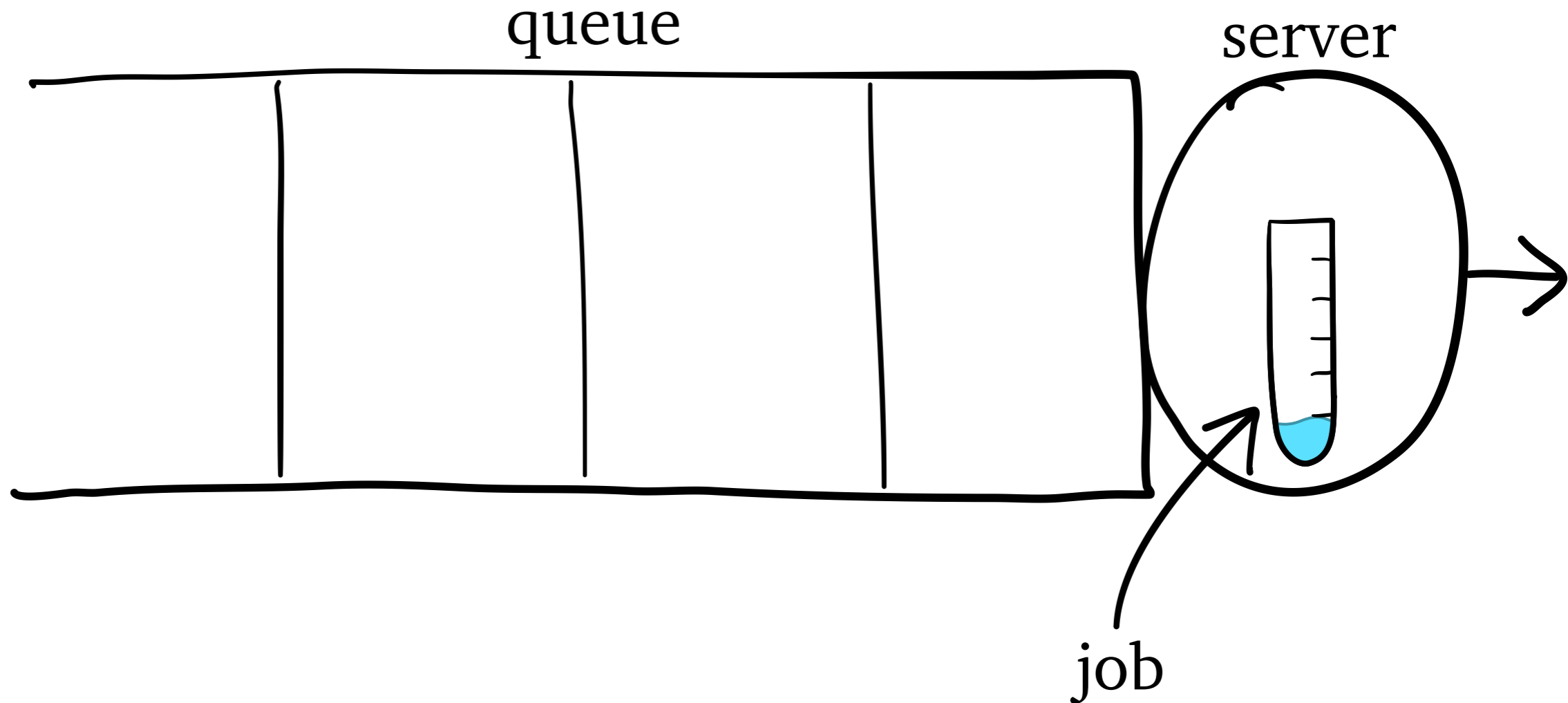
# M/G/1 Queue



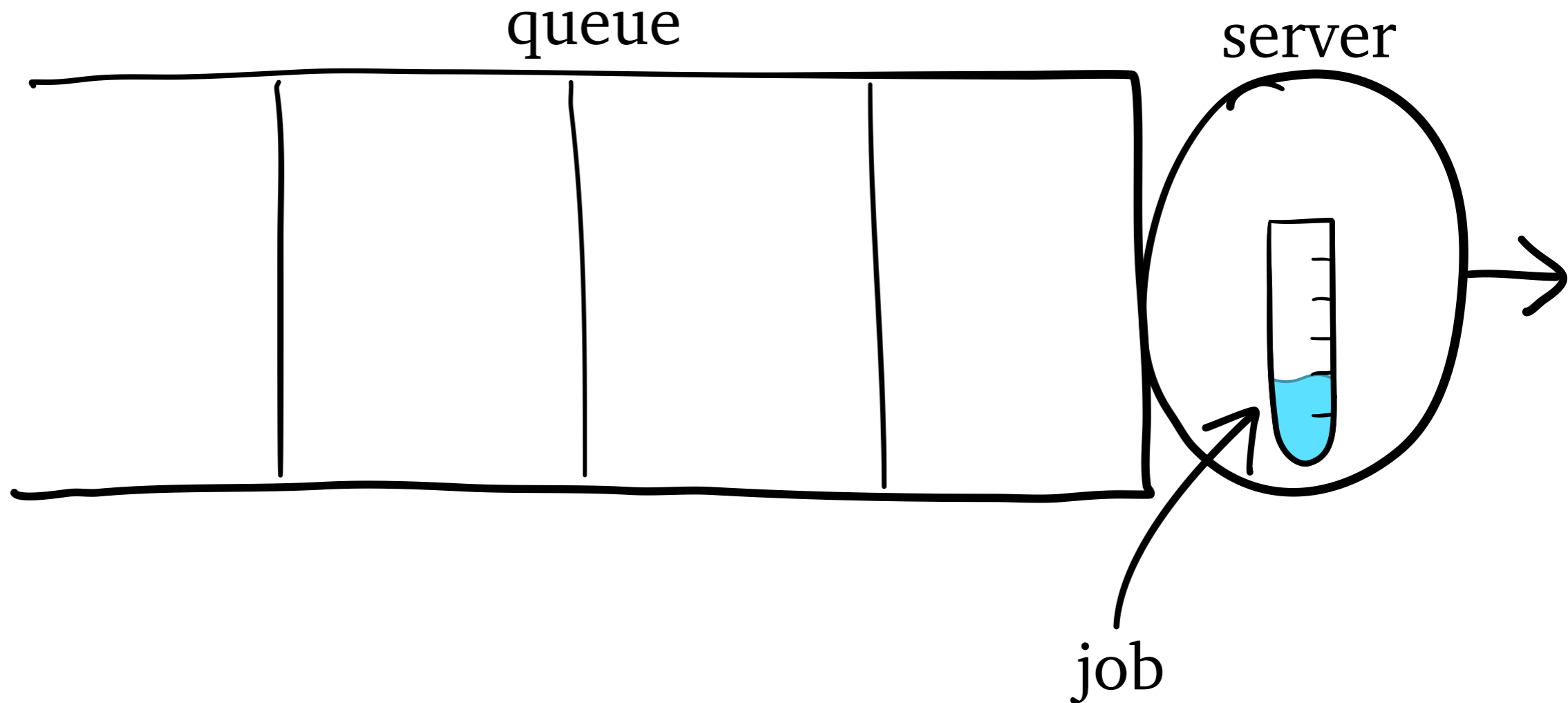
# M/G/1 Queue



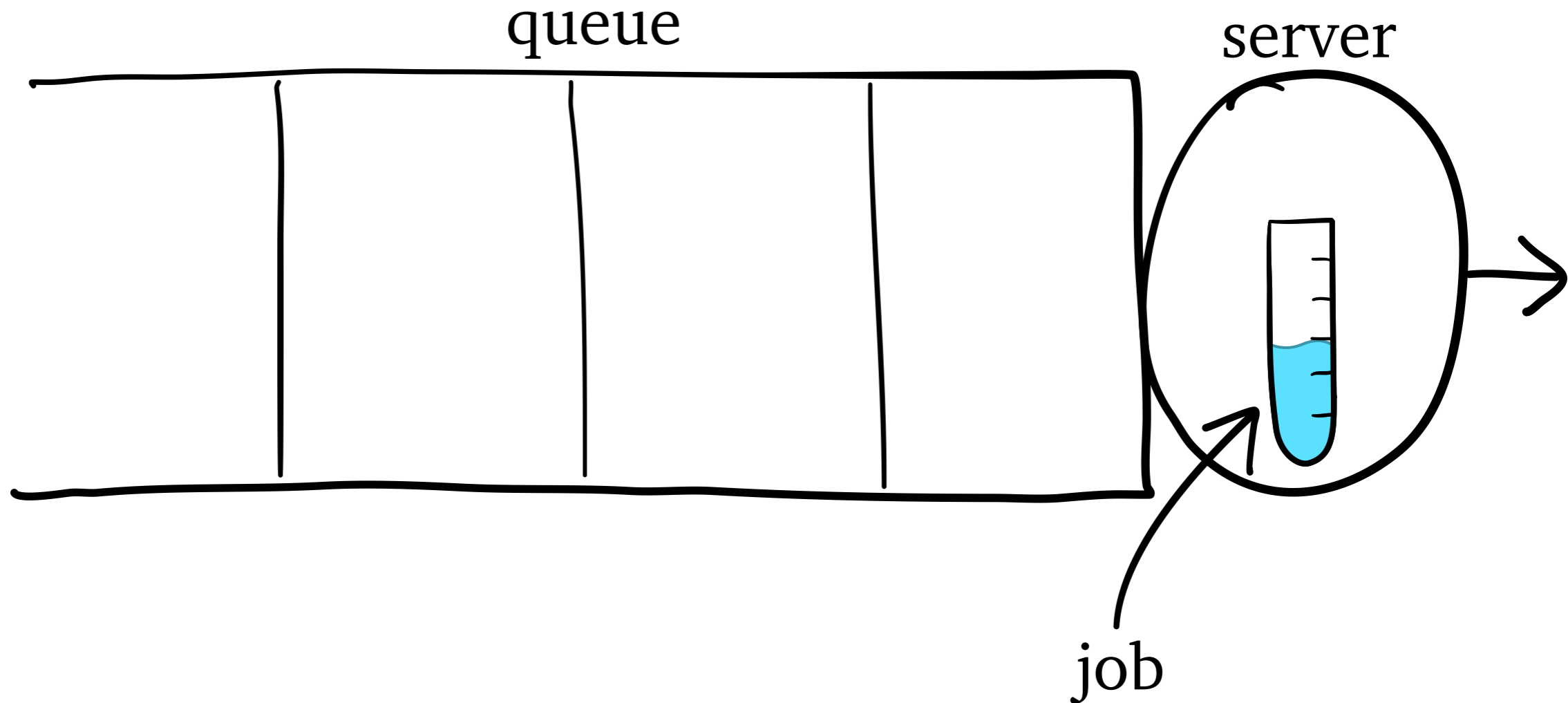
# M/G/1 Queue



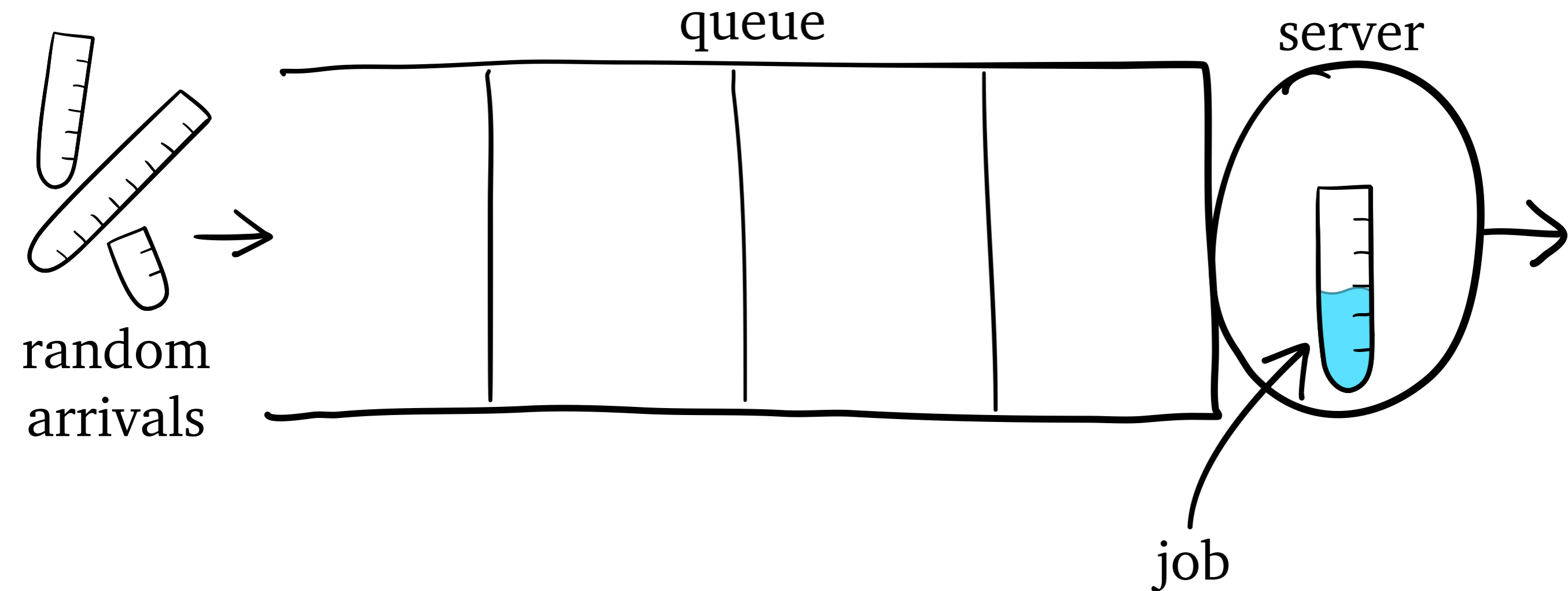
# M/G/1 Queue



# M/G/1 Queue

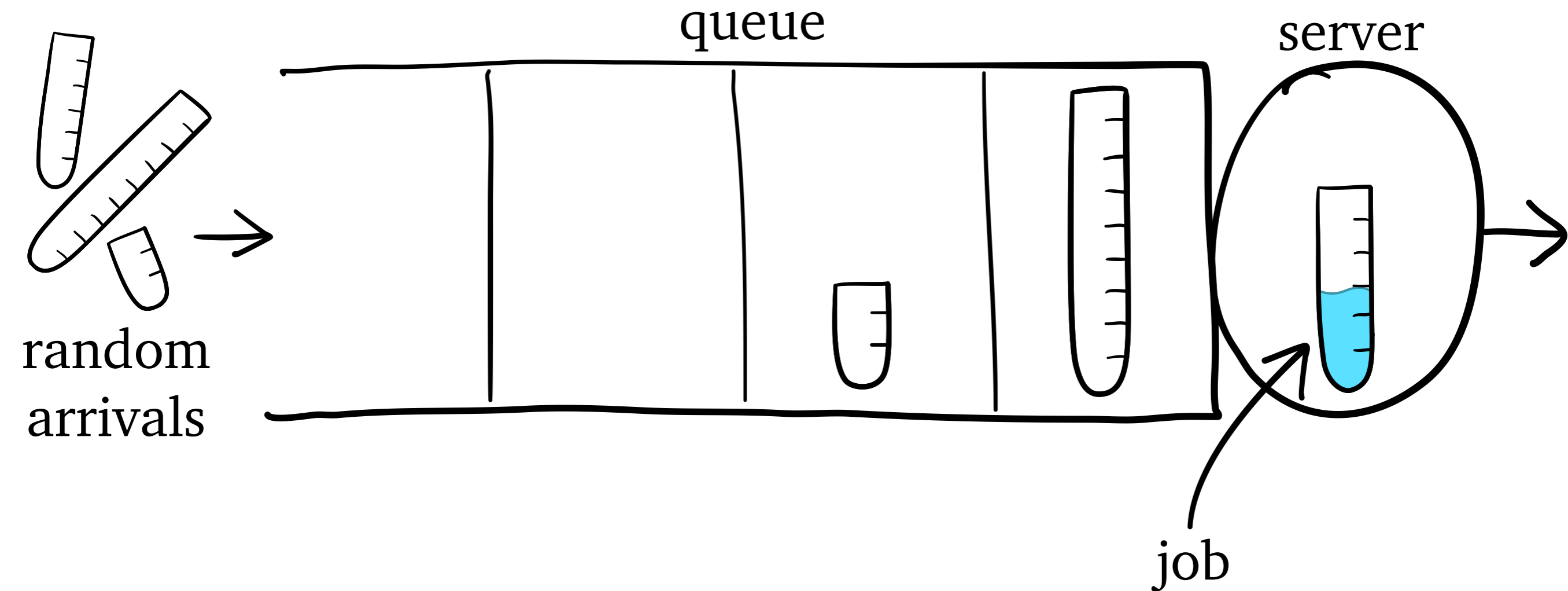


# M/G/1 Queue

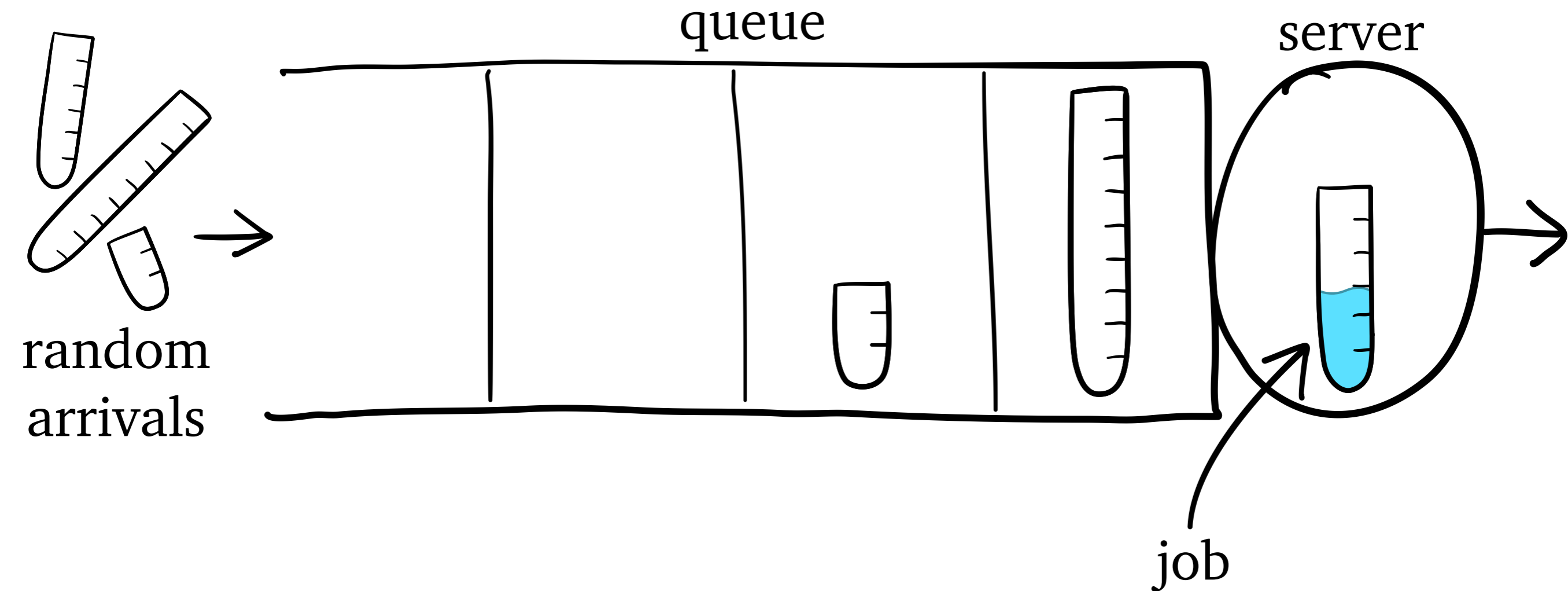




# M/G/1 Queue

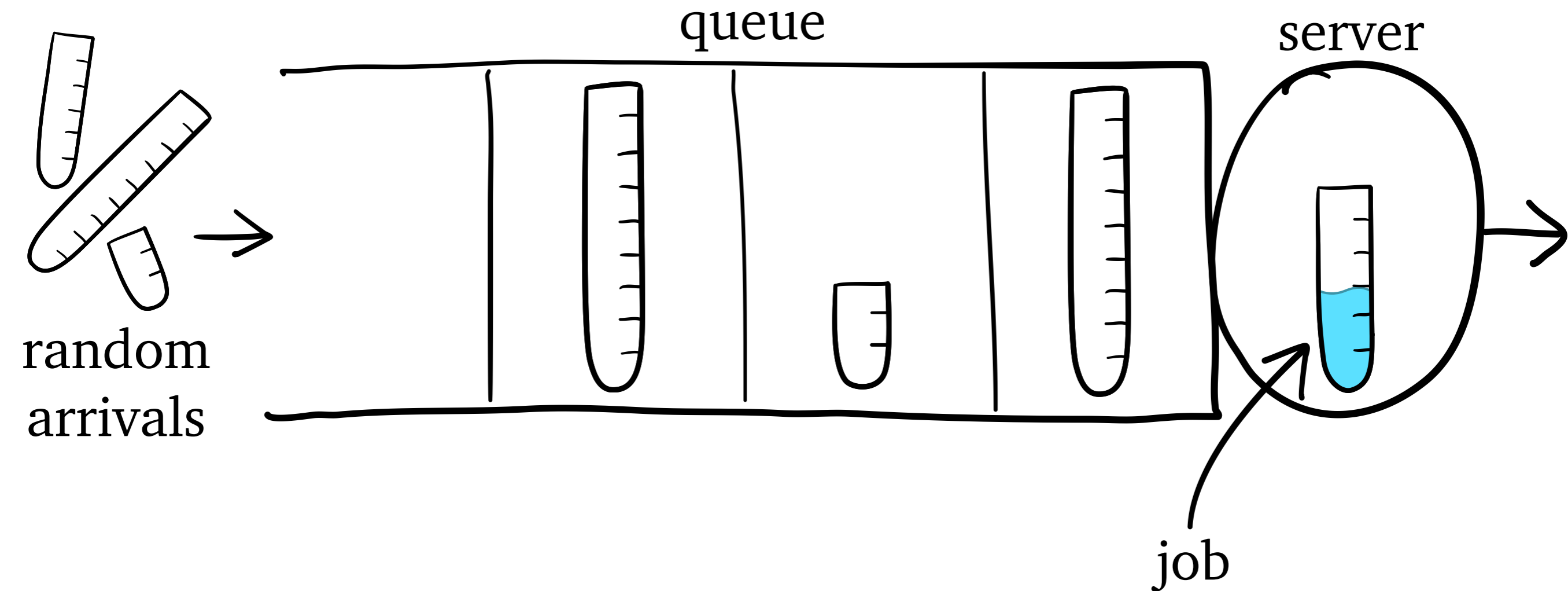


# M/G/1 Queue



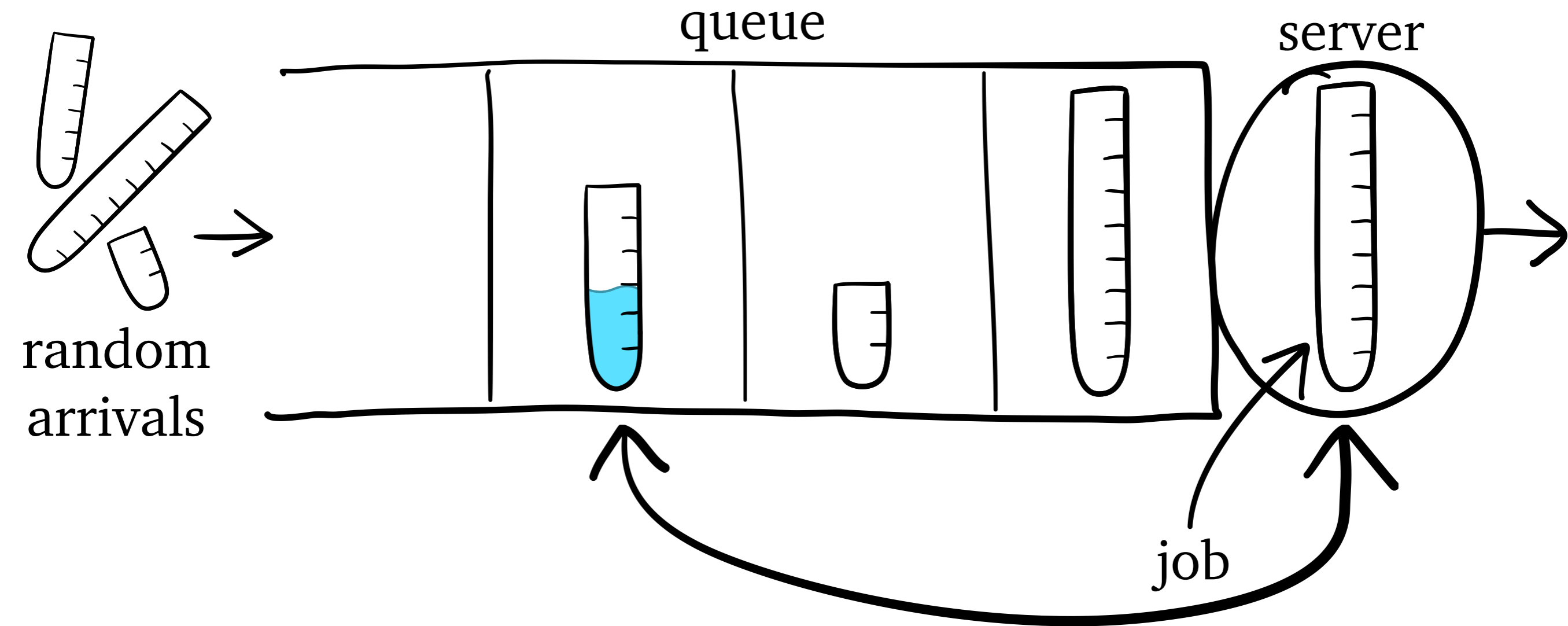
**Scheduling policy:** picks which job to serve

# M/G/1 Queue



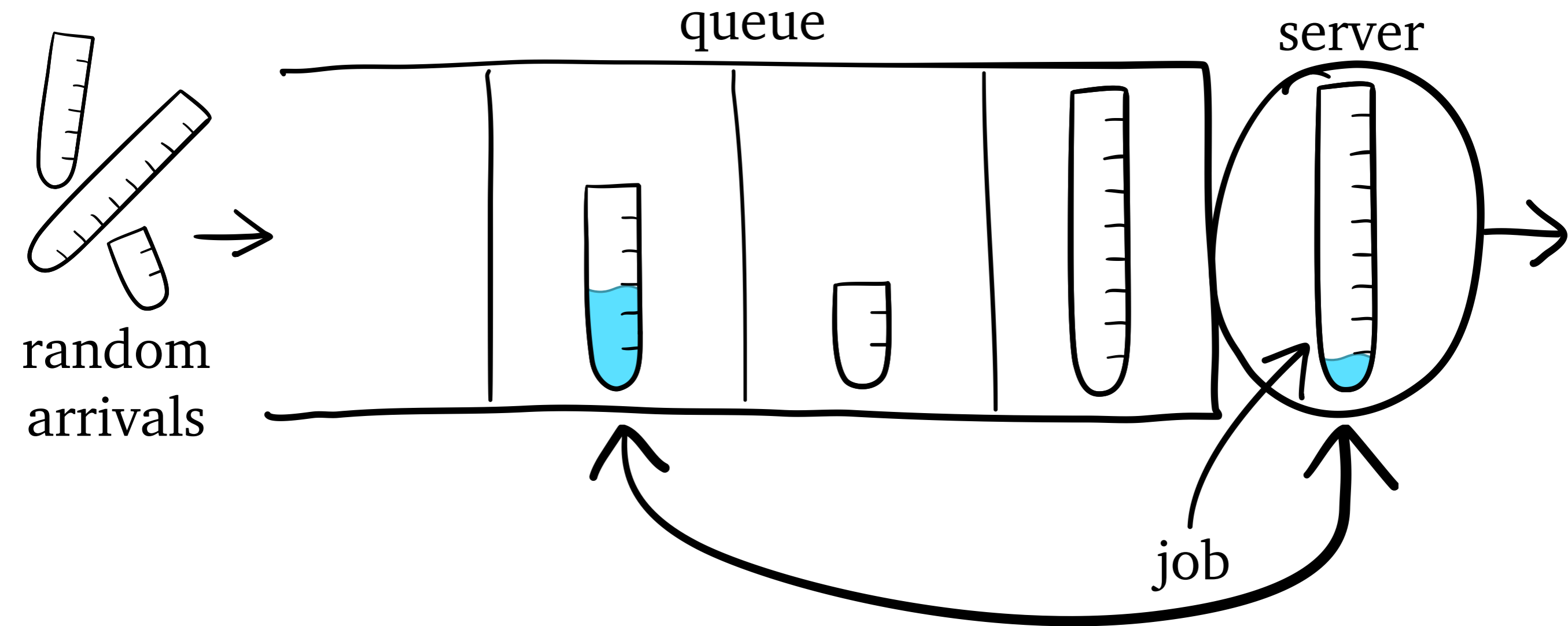
**Scheduling policy:** picks which job to serve

# M/G/1 Queue



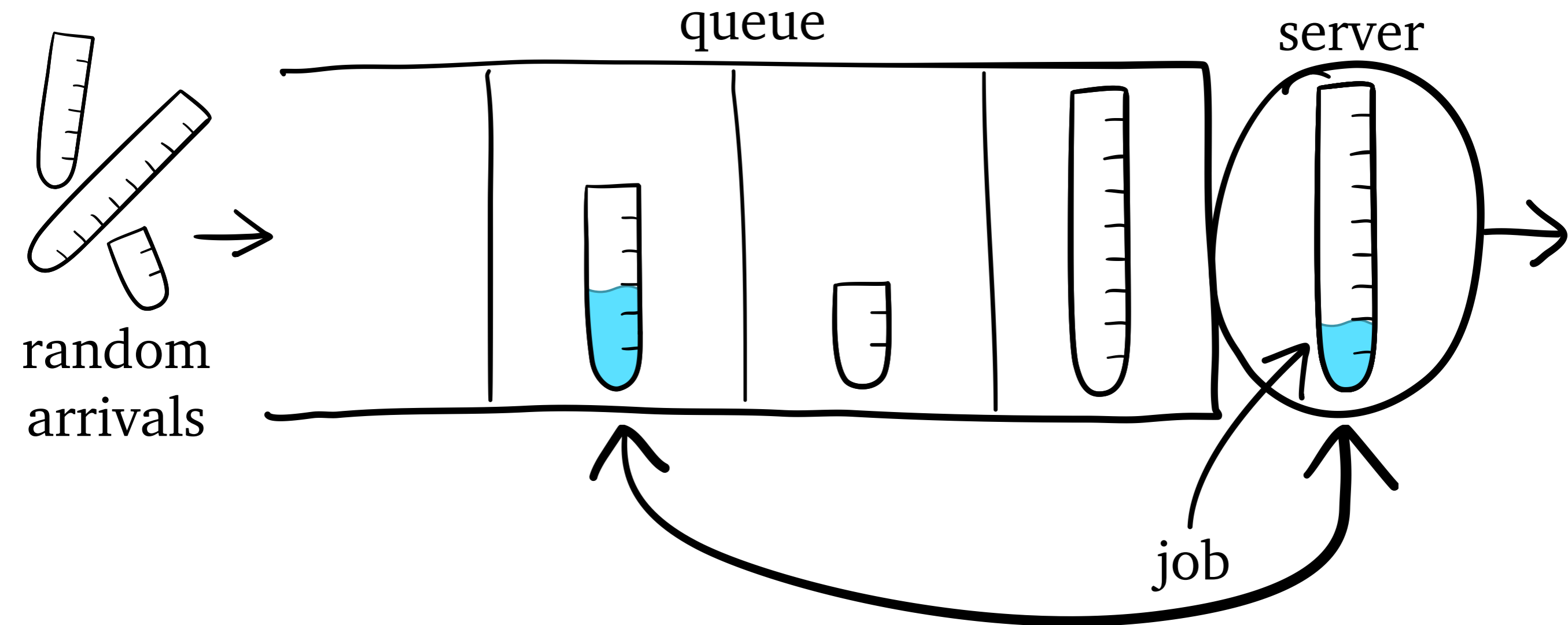
**Scheduling policy:** picks which job to serve

# M/G/1 Queue



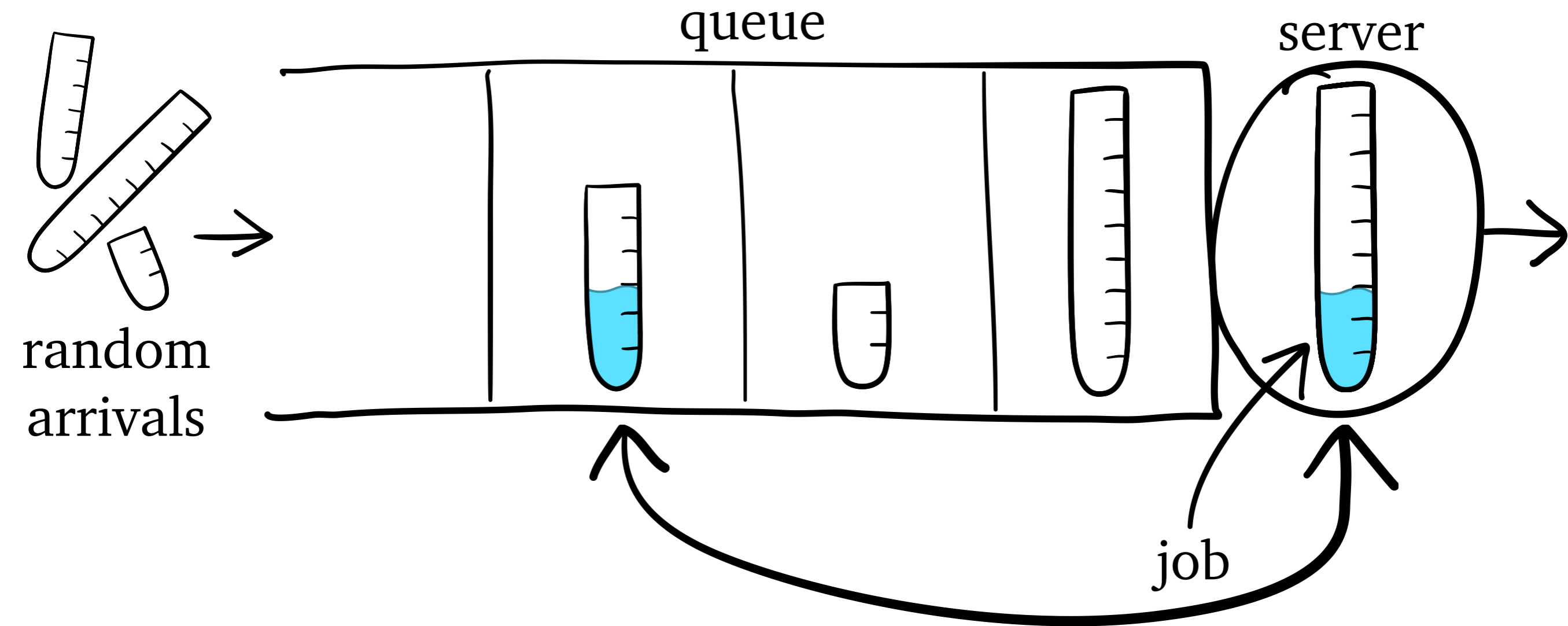
**Scheduling policy:** picks which job to serve

# M/G/1 Queue



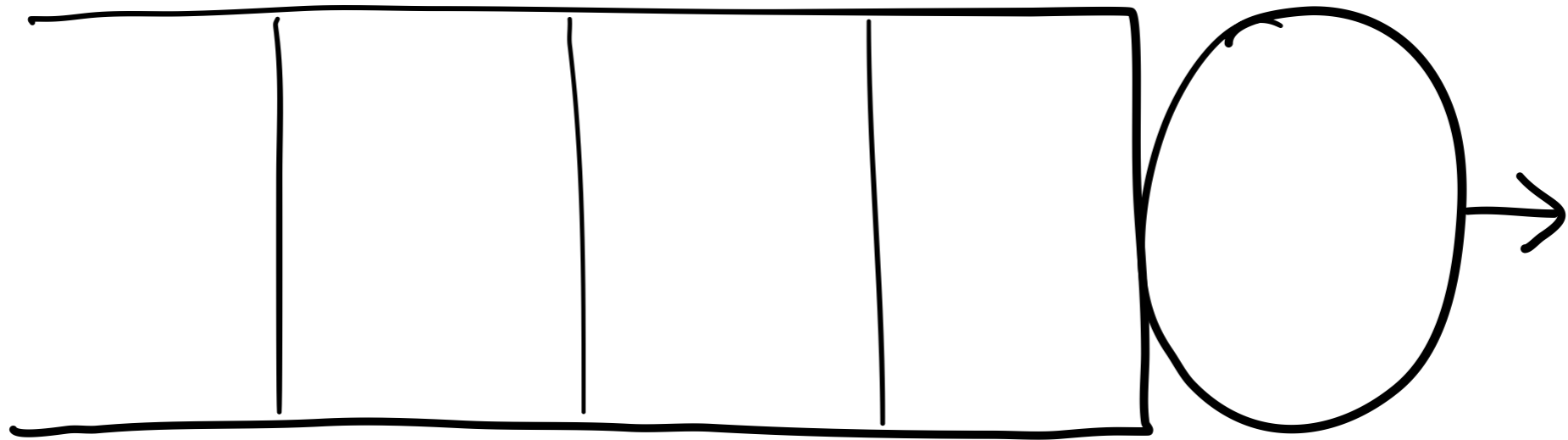
**Scheduling policy:** picks which job to serve

# M/G/1 Queue



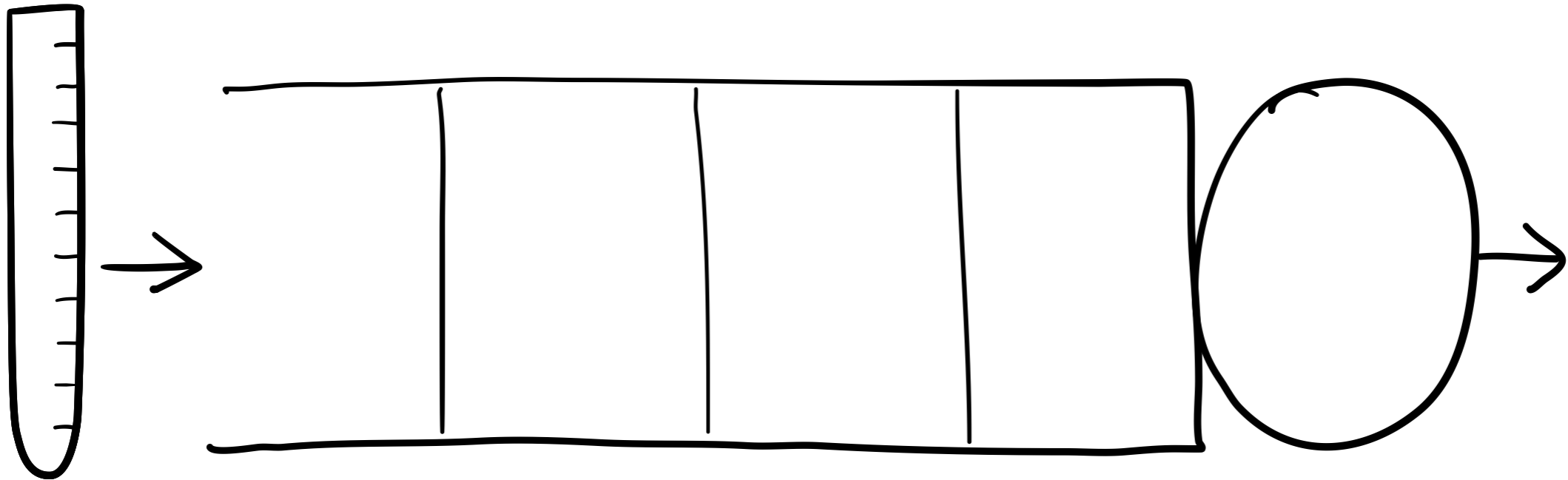
**Scheduling policy:** picks which job to serve

# Response Time

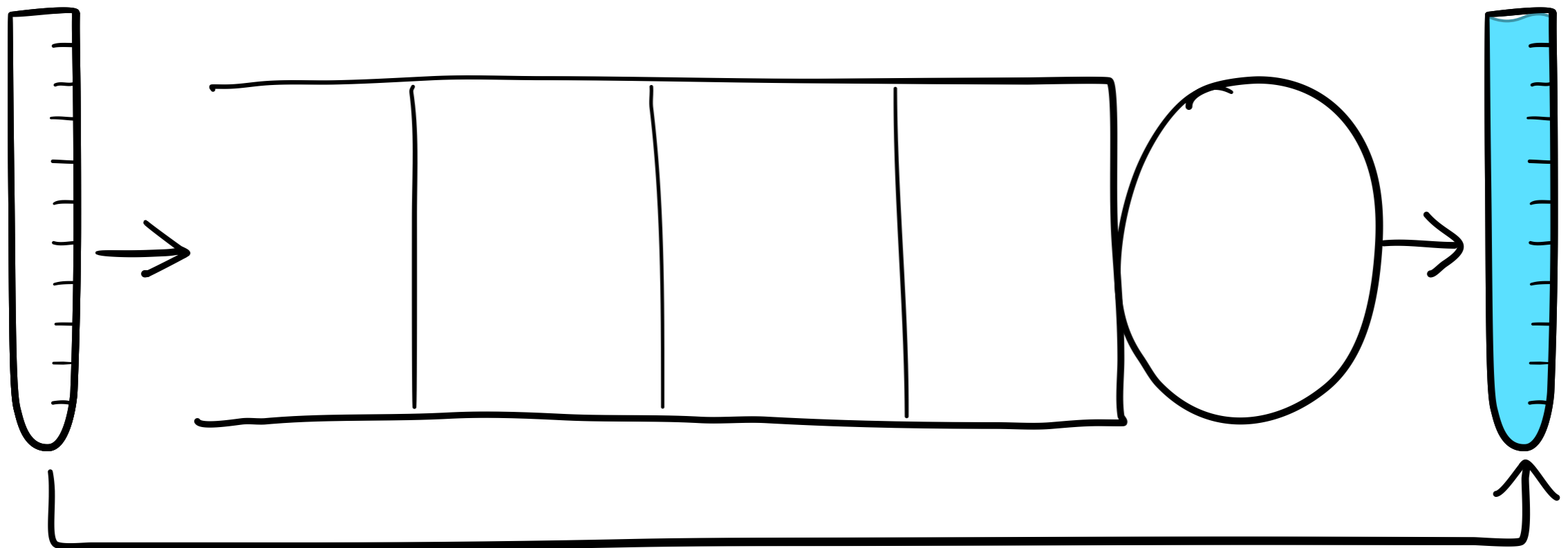




# Response Time

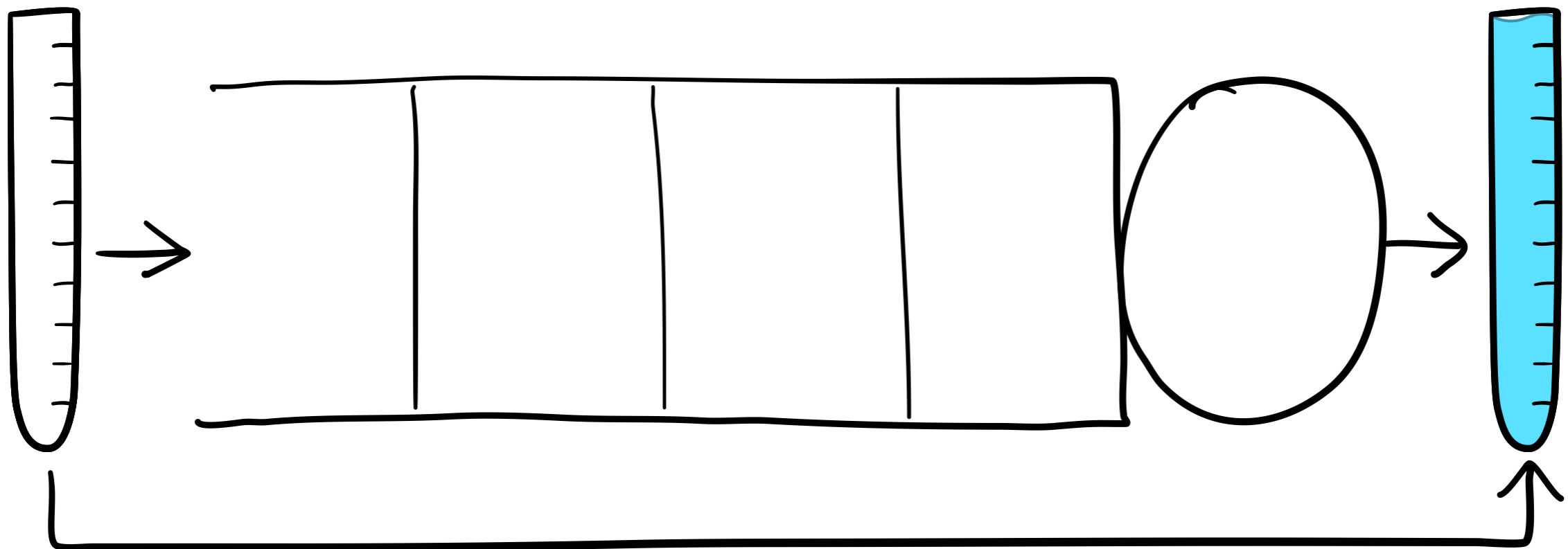


# Response Time



$T = \text{response time}$

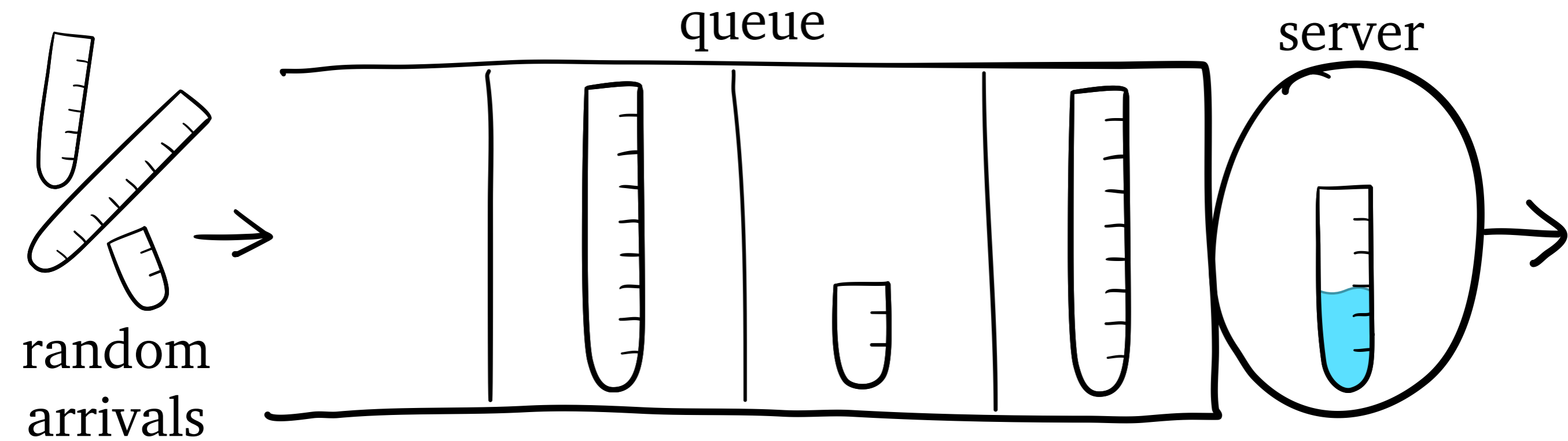
# Response Time



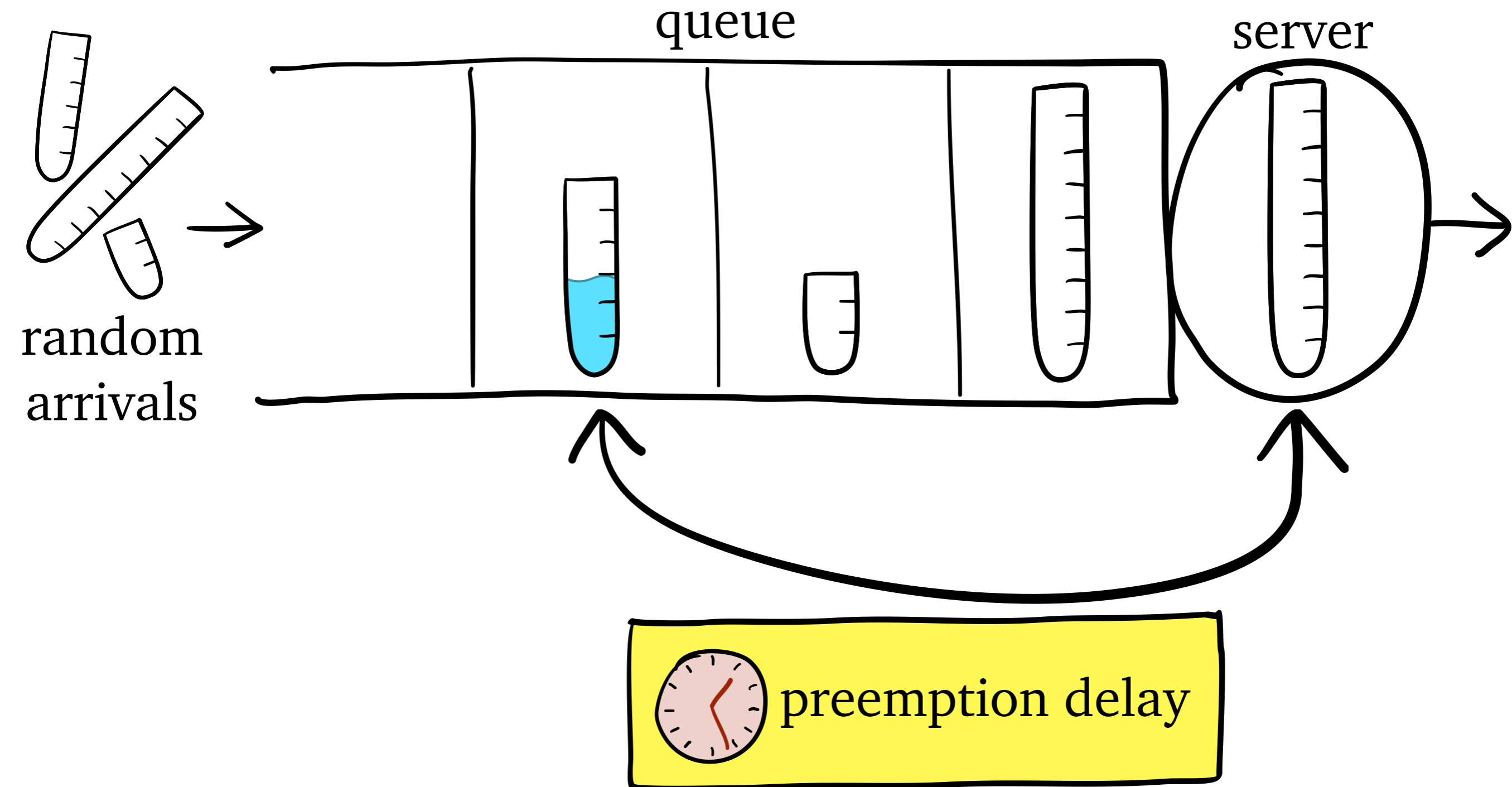
$T = \text{response time}$

**Goal:** analyze and minimize *mean response time*  $\mathbf{E}[T]$

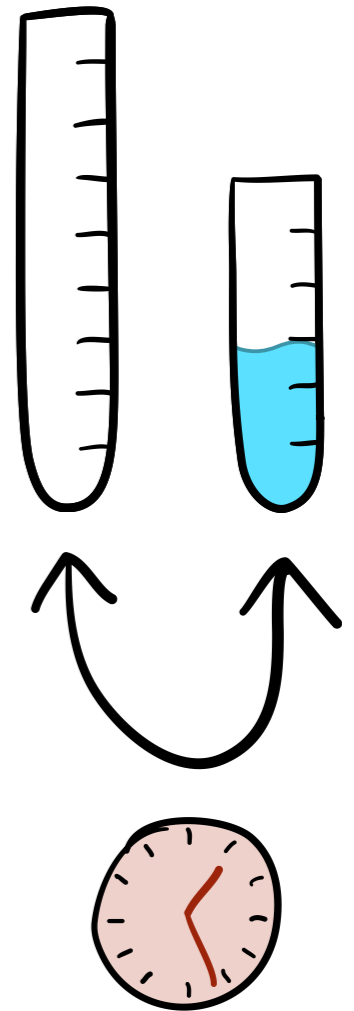
# M/G/1 Queue



# M/G/1 Queue

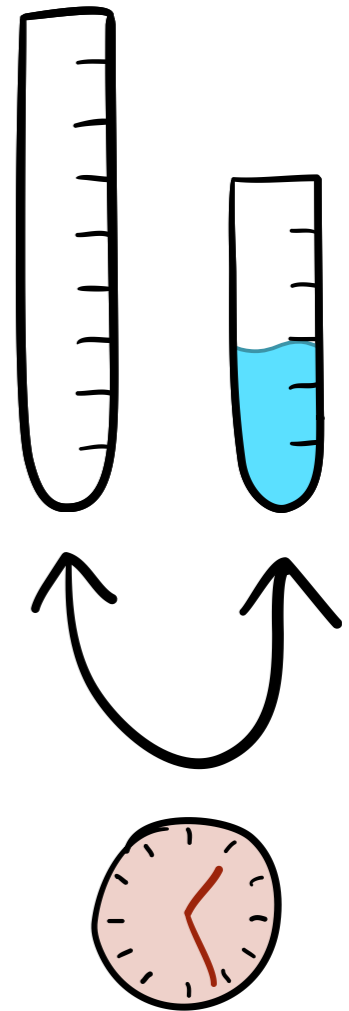


# What is $E[T]$ with Delays?



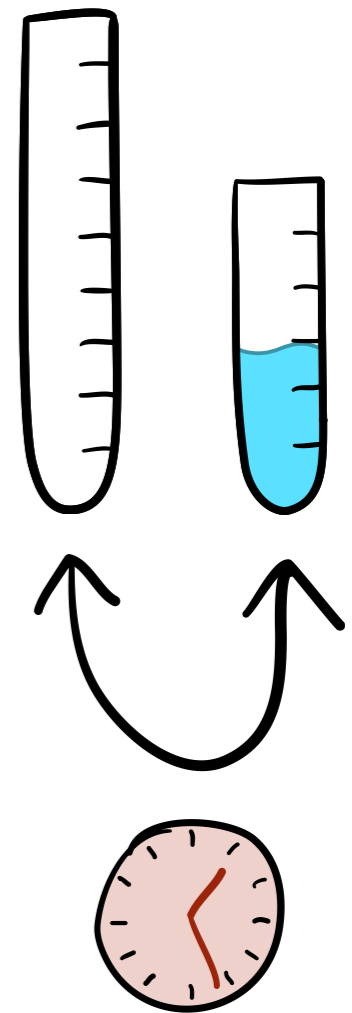
# What is $E[T]$ with Delays?

- Any nonpreemptive policy:  
add delay to job size



# What is $E[T]$ with Delays?

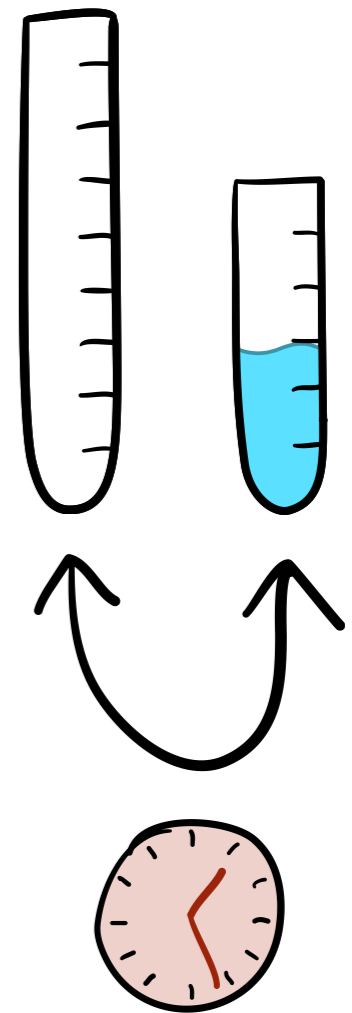
- Any nonpreemptive policy:  
add delay to job size
- Any preemptive policy:  
even *stability* hard to know



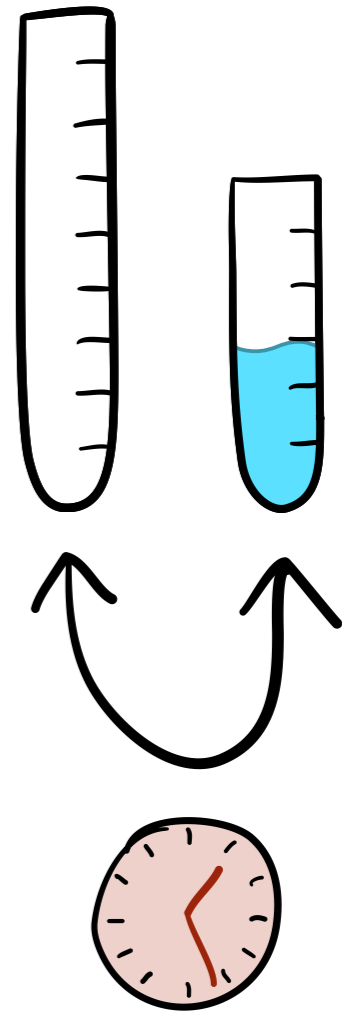


# What is $E[T]$ with Delays?

- Any nonpreemptive policy:  
add delay to job size
- Any preemptive policy:  
even *stability* hard to know
- Only known analysis: SRPT [Goerg '86]

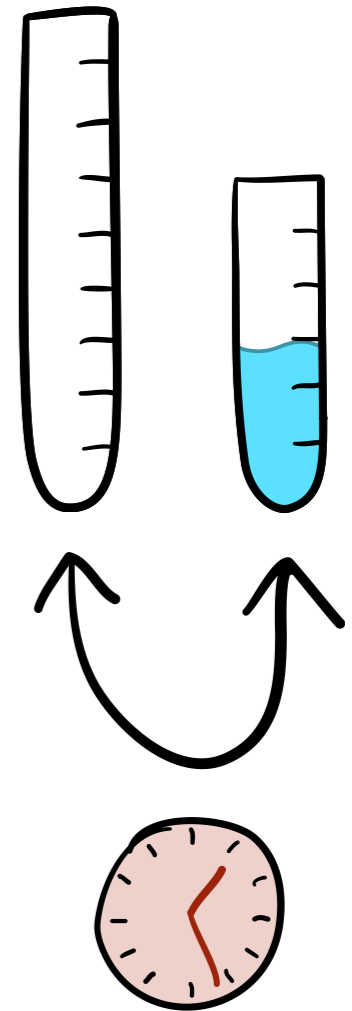


# Is There a Stable Policy?



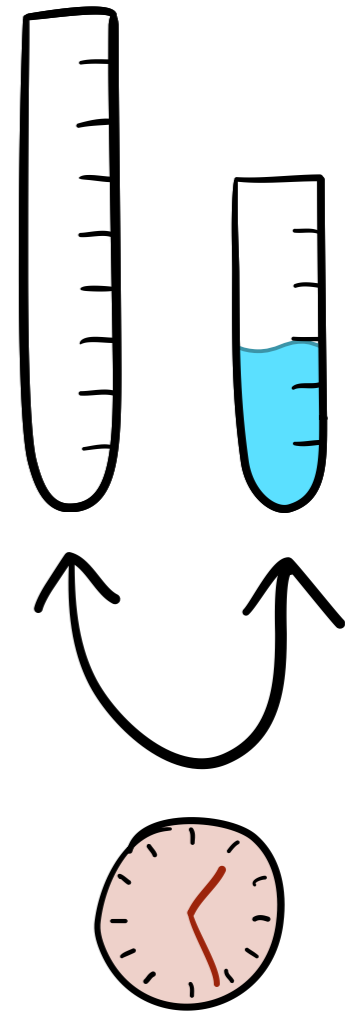
# Is There a Stable Policy?

- When load is near 1, *preemptive* policies can raise effective load to over 1



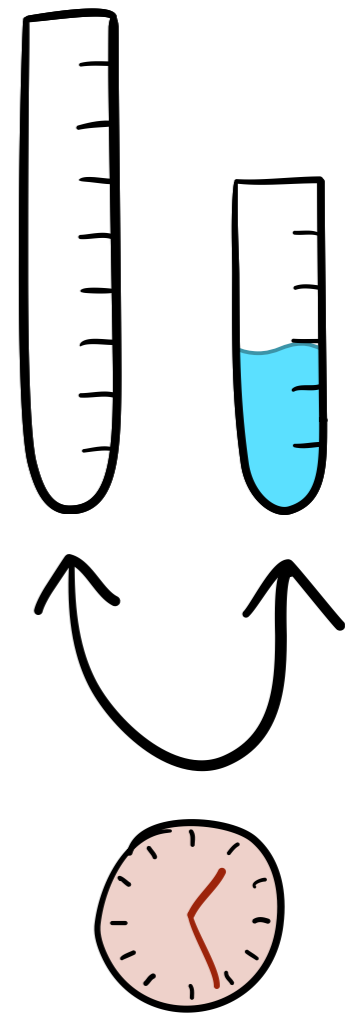
# Is There a Stable Policy?

- When load is near 1, *preemptive* policies can raise effective load to over 1
  - *load instability*



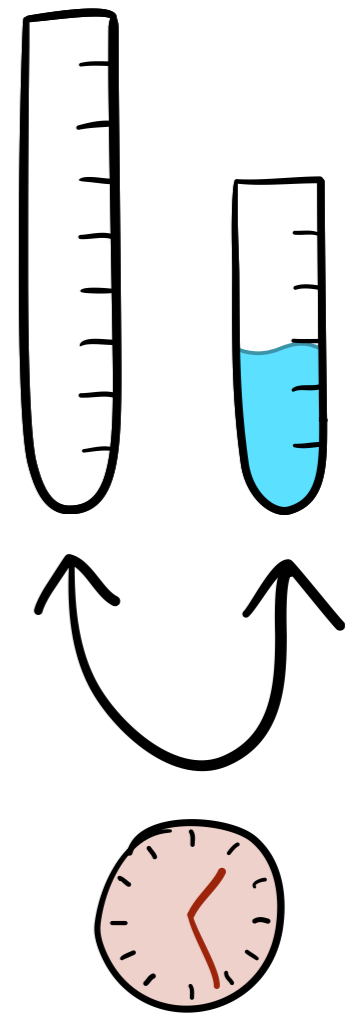
# Is There a Stable Policy?

- When load is near 1, *preemptive* policies can raise effective load to over 1
  - *load instability*
- When  $E[X^2]$  is infinite, *nonpreemptive* policies have infinite  $E[T]$



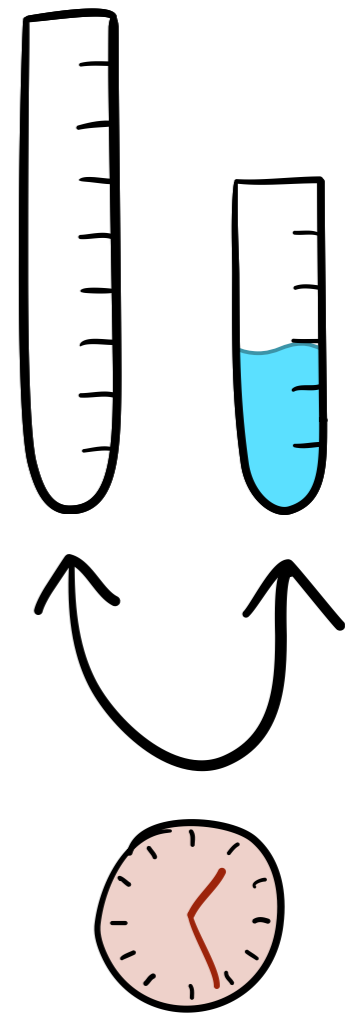
# Is There a Stable Policy?

- When load is near 1, *preemptive* policies can raise effective load to over 1
  - *load instability*
- When  $E[X^2]$  is infinite, *nonpreemptive* policies have infinite  $E[T]$ 
  - *response time instability*



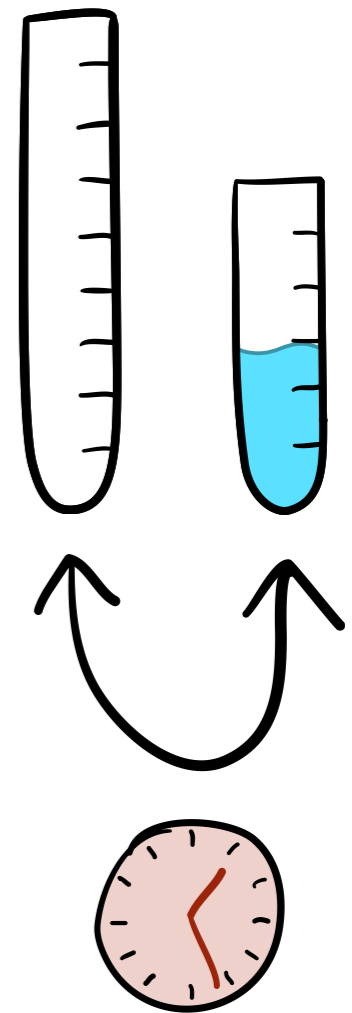
# Is There a Stable Policy?

- When load is near 1, *preemptive* policies can raise effective load to over 1
  - *load instability*
- When  $E[X^2]$  is infinite, *nonpreemptive* policies have infinite  $E[T]$ 
  - *response time instability*
- Is there a policy that has both?



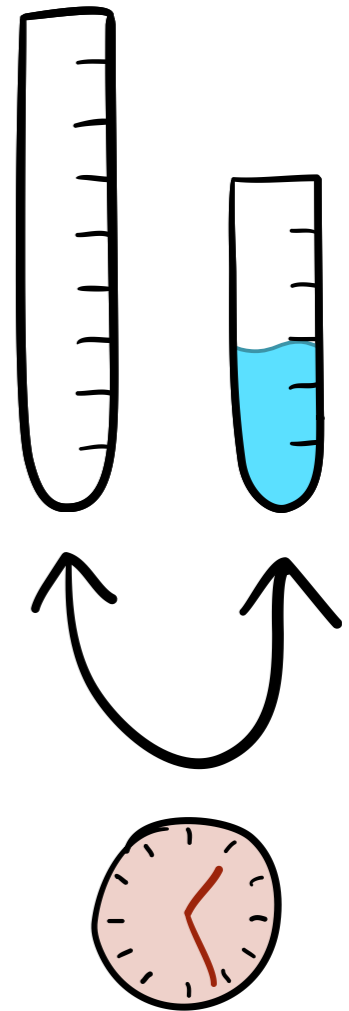
# Is There a Stable Policy?

- When load is near 1, *preemptive* policies can raise effective load to over 1
  - *load instability*
- When  $E[X^2]$  is infinite, *nonpreemptive* policies have infinite  $E[T]$ 
  - *response time instability*
- Is there a policy that has both?
  - Probably a variant of PLCFS?



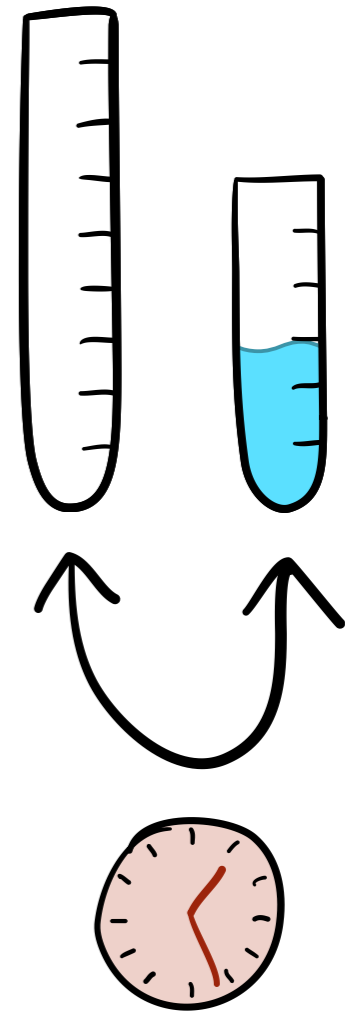


# How do we Minimize $E[T]$ ?



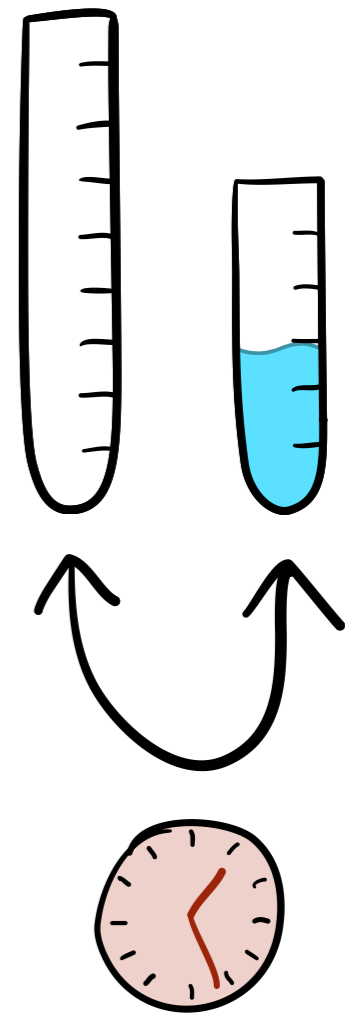
# How do we Minimize $E[T]$ ?

- Policy for no preemption delay: *Gittins policy*



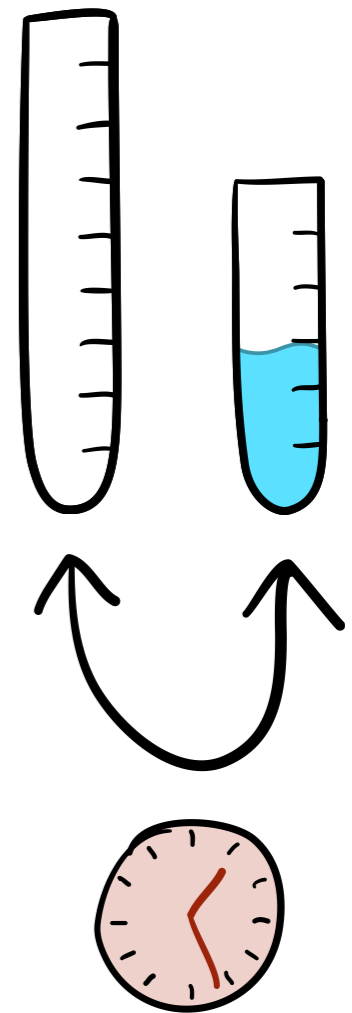
# How do we Minimize $E[T]$ ?

- Policy for no preemption delay: *Gittins policy*
  - Comes from multiarmed bandits



# How do we Minimize $E[T]$ ?

- Policy for no preemption delay: *Gittins policy*
  - Comes from multiarmed bandits
- Probably related: multiarmed bandits with switching costs [Asawa and Teneketzis '96]



# How do we Minimize $E[T]$ ?

- Policy for no preemption delay: *Gittins policy*
  - Comes from multiarmed bandits
- Probably related: multiarmed bandits with swicthing costs [Asawa and Teneketzis '96]
  - ... but *cost* and *delay* are different

