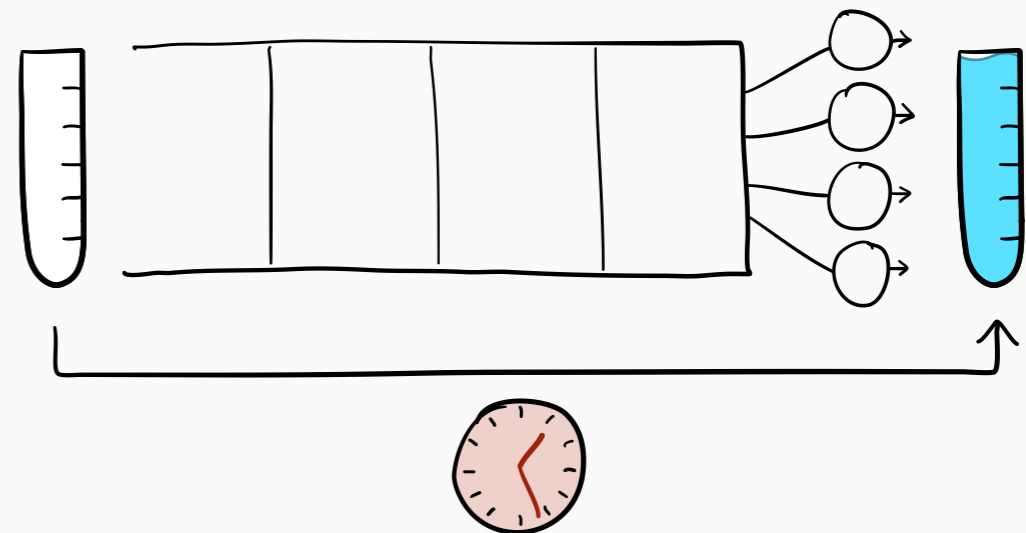# The Gittins Policy Is Nearly Optimal in the M/G/*k*

*under Extremely General Conditions*
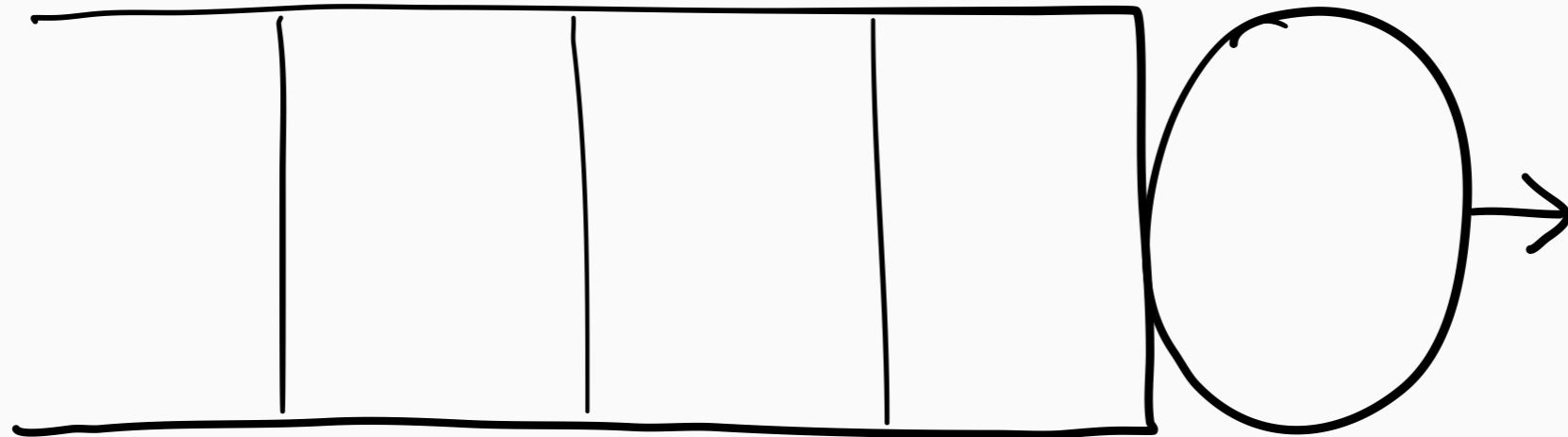
Ziv Scully
Isaac Grosof
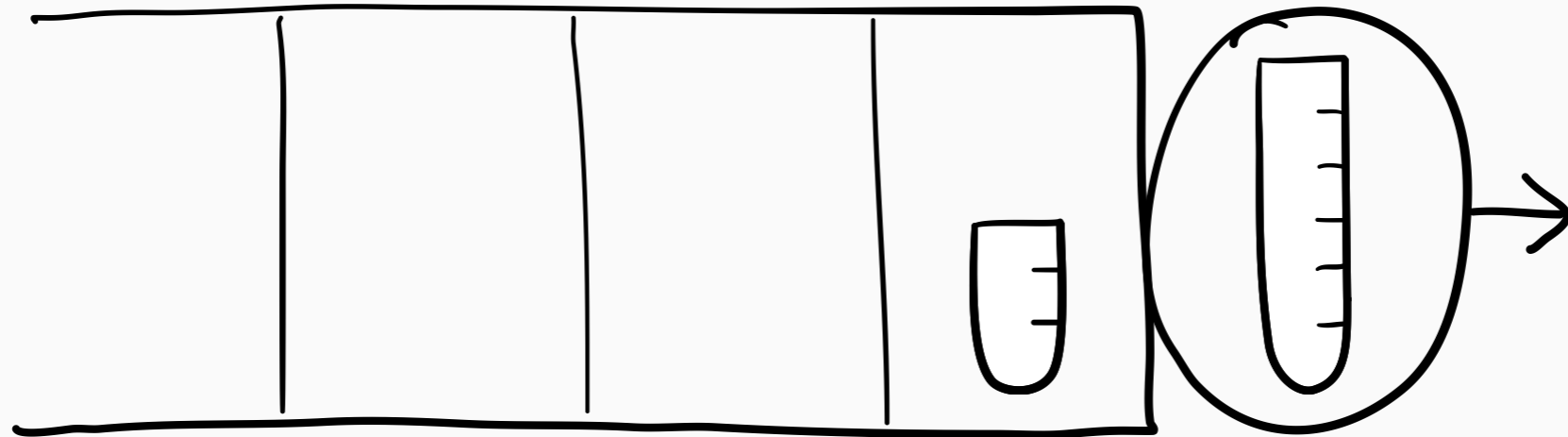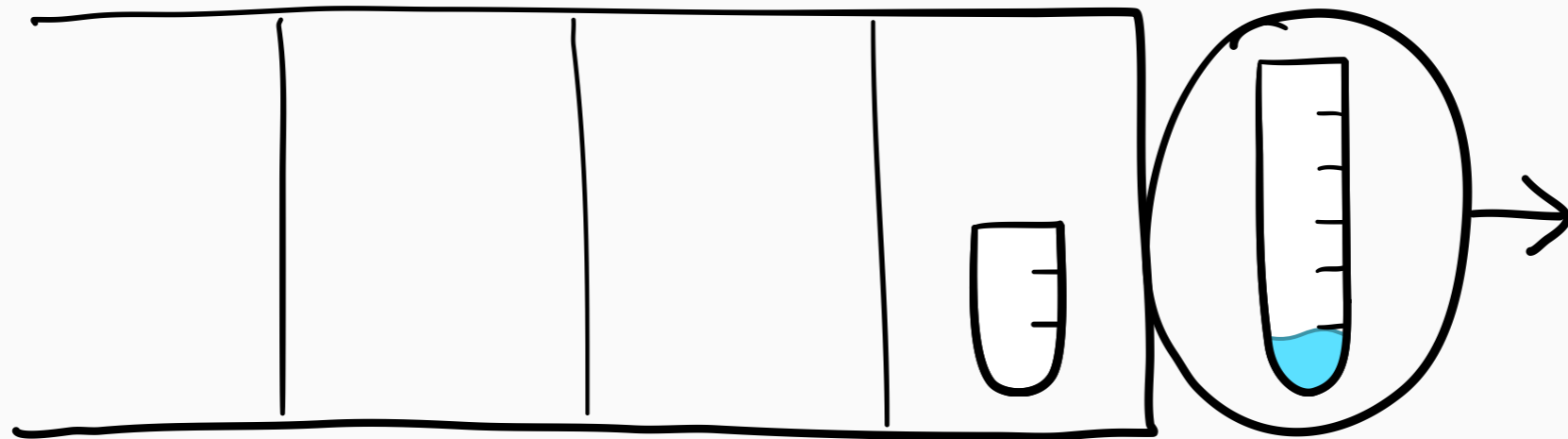Mor Harchol-Balter

*Carnegie Mellon University*

# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?

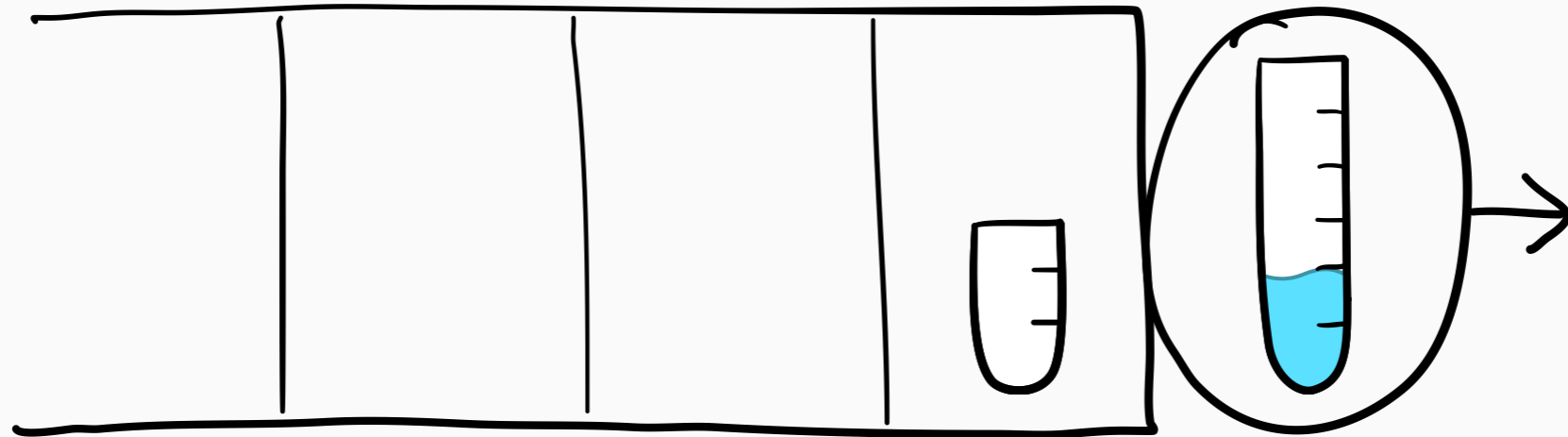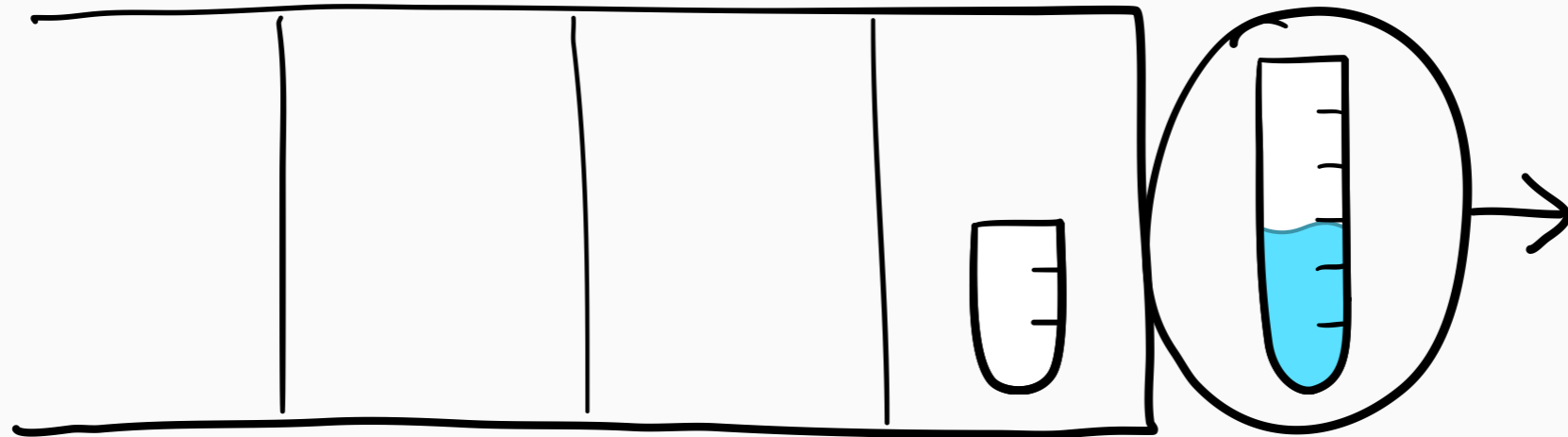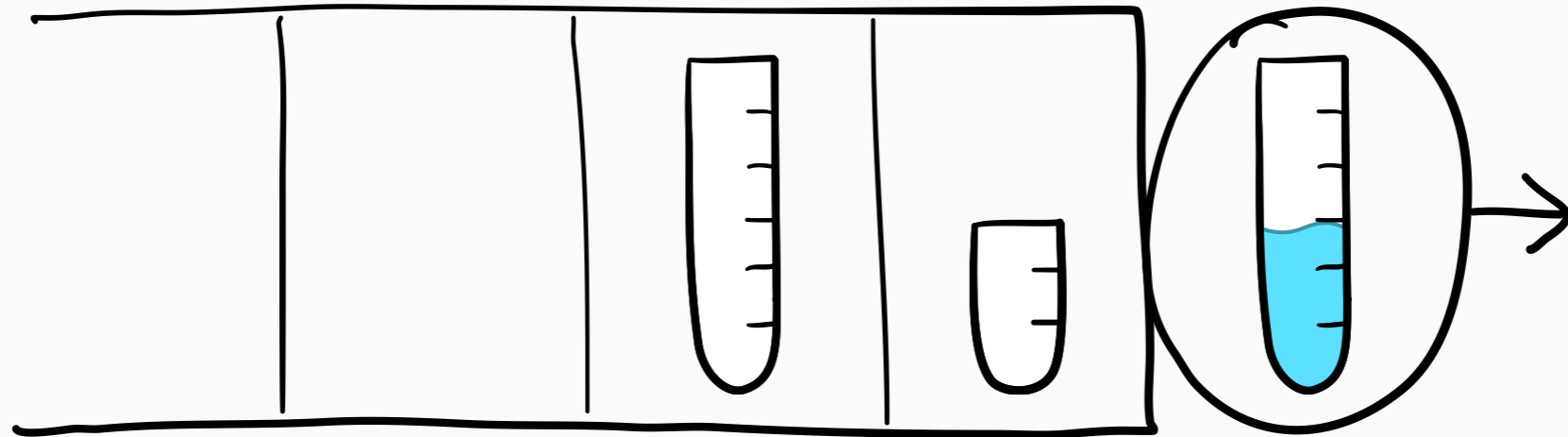# How should we schedule jobs
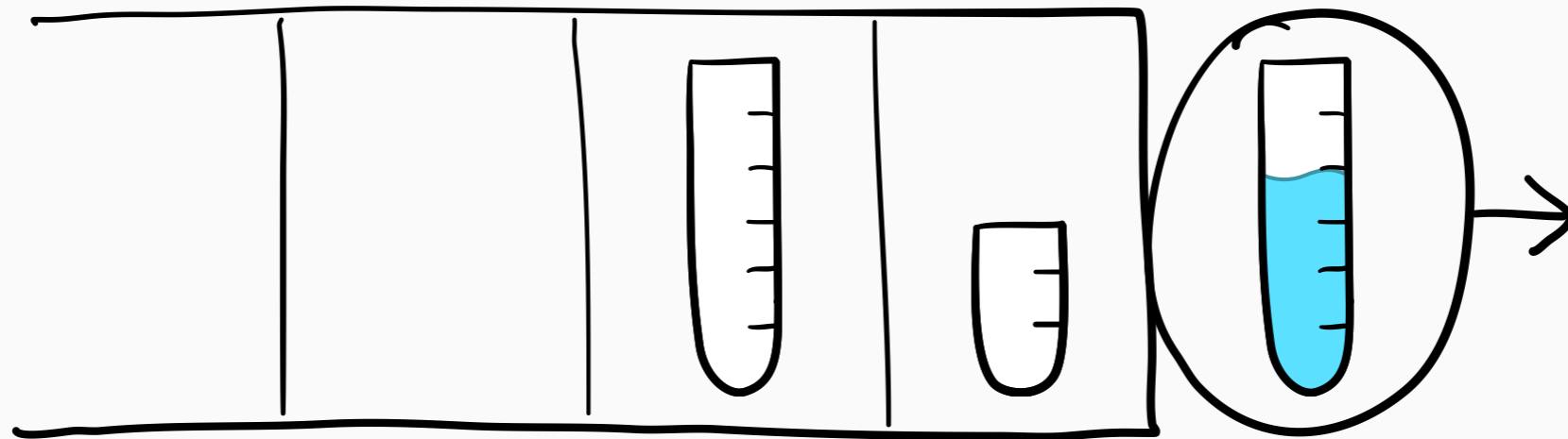# to minimize delay?

# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?
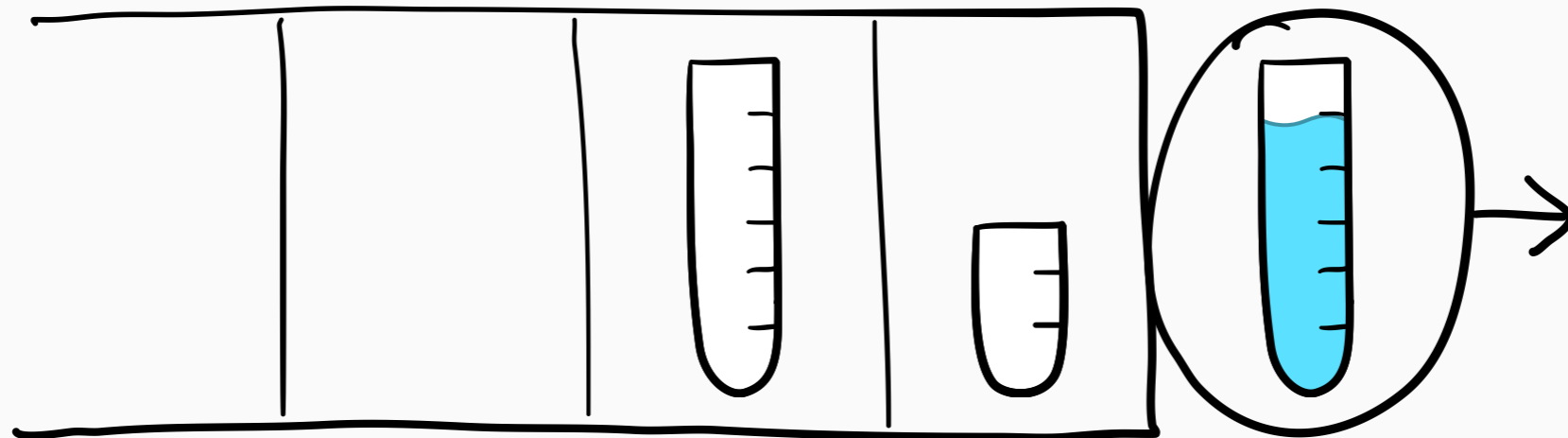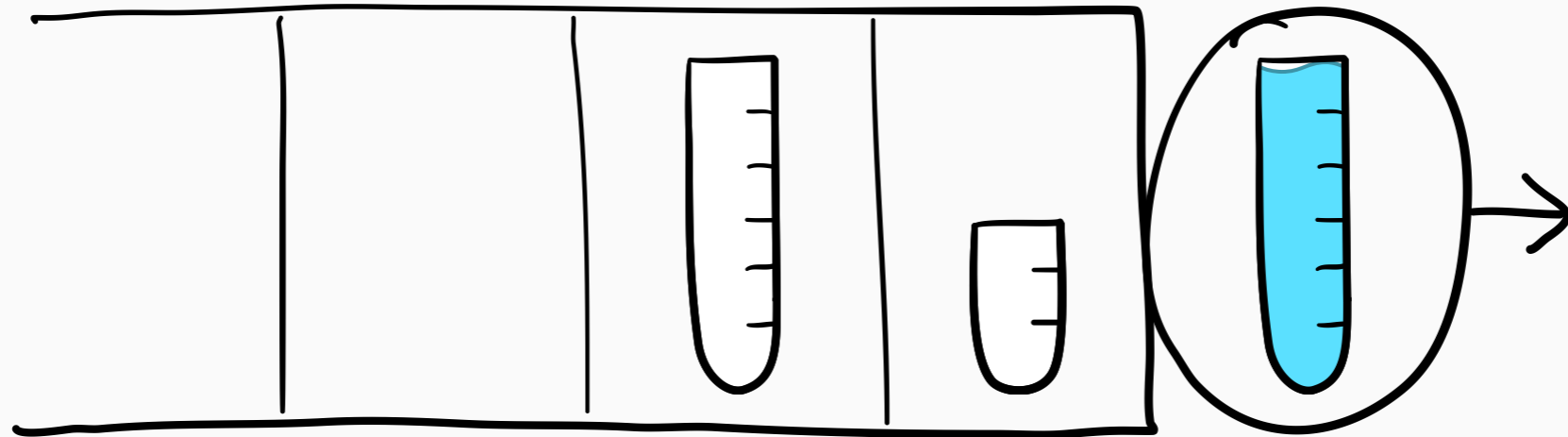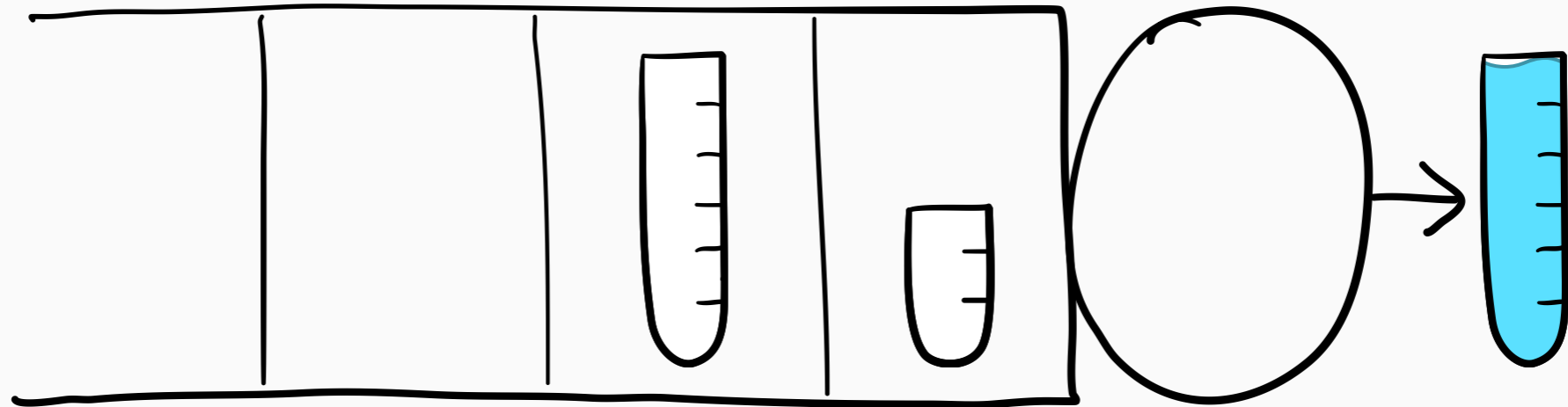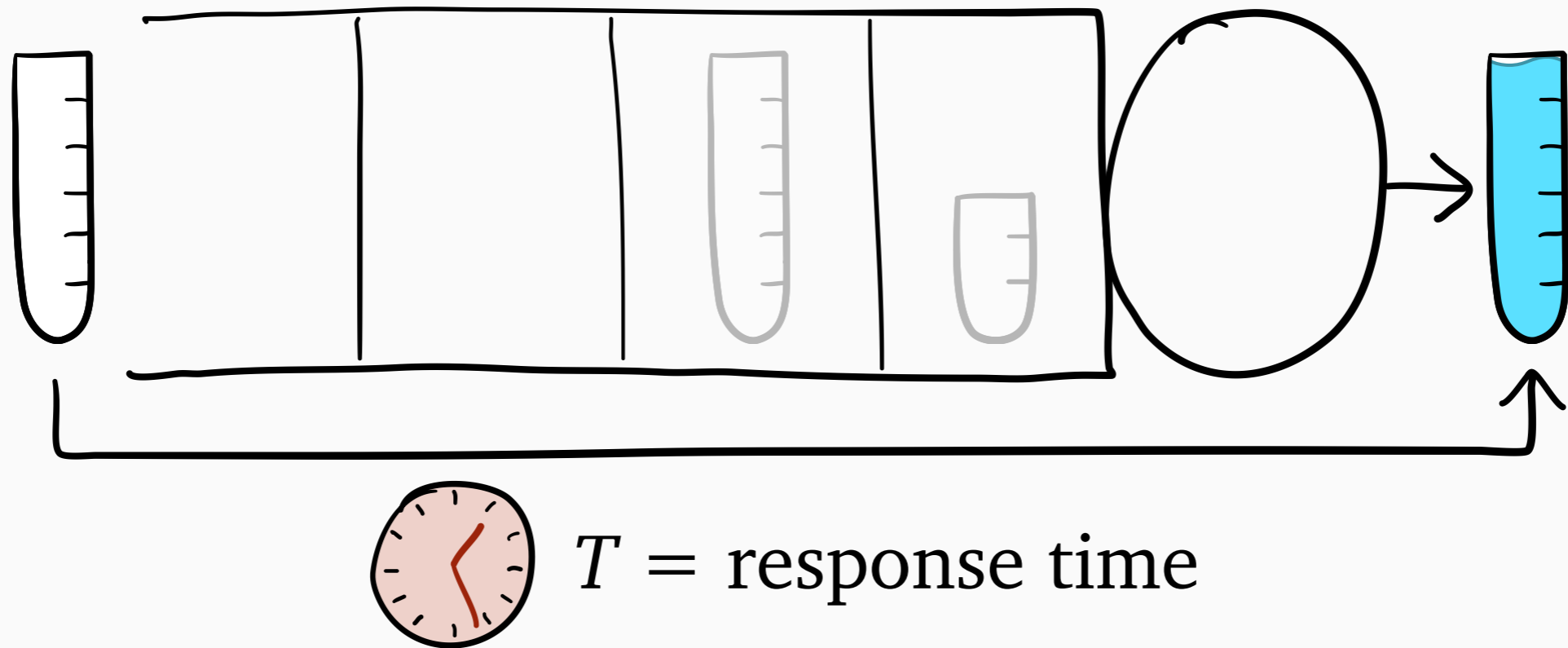
# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?

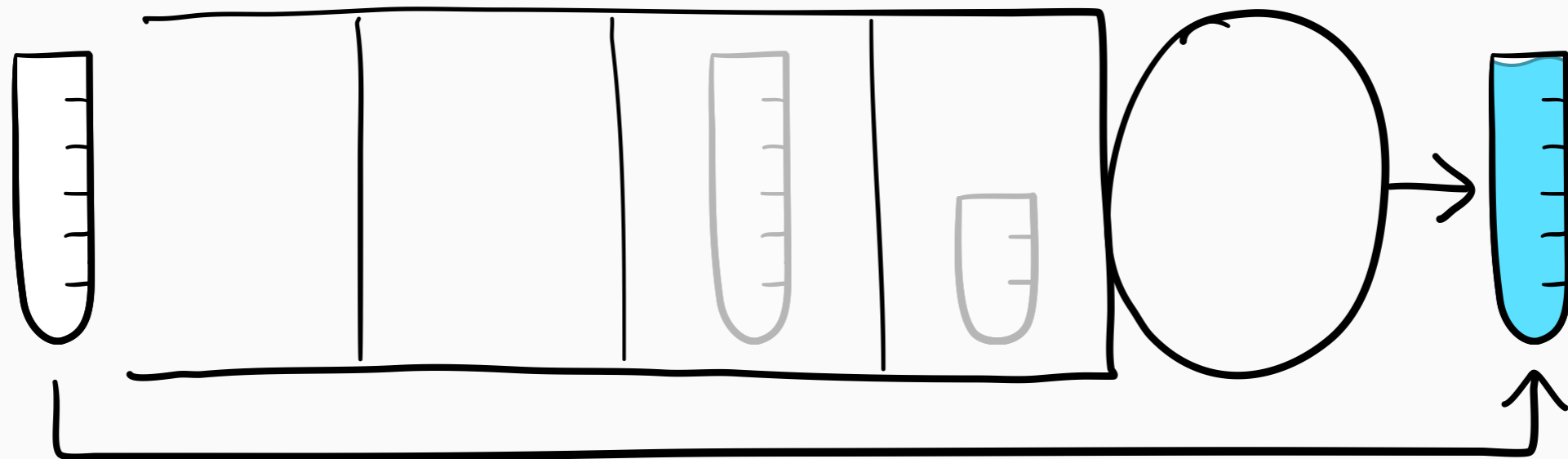# How should we schedule jobs to minimize delay?

# How should we schedule jobs to minimize delay?
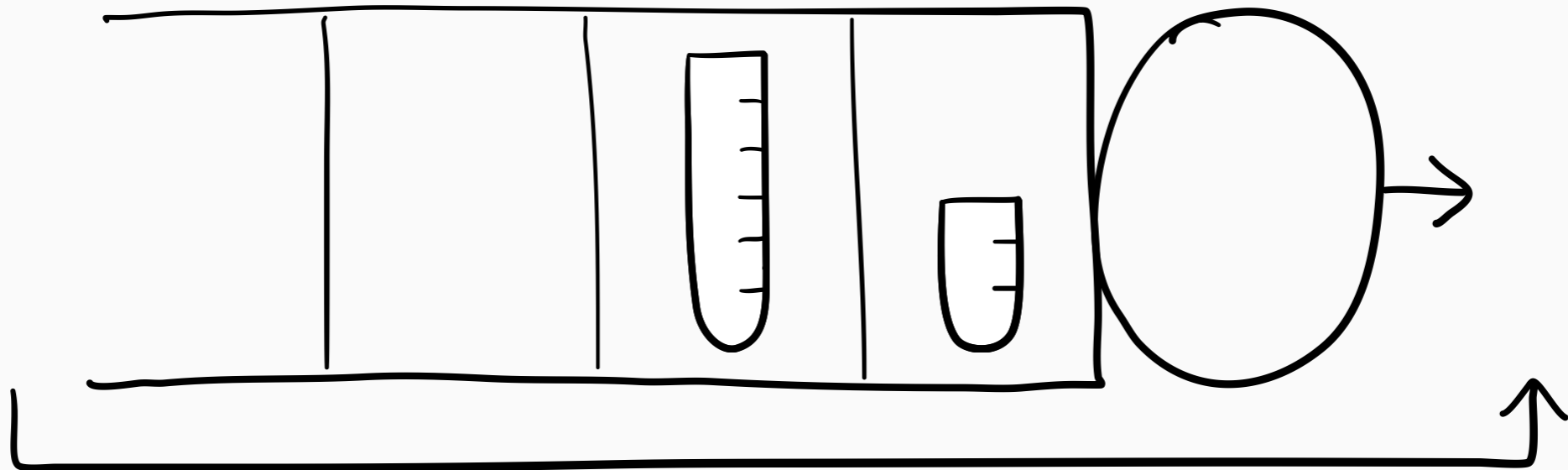


$T$ = response time
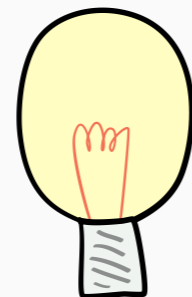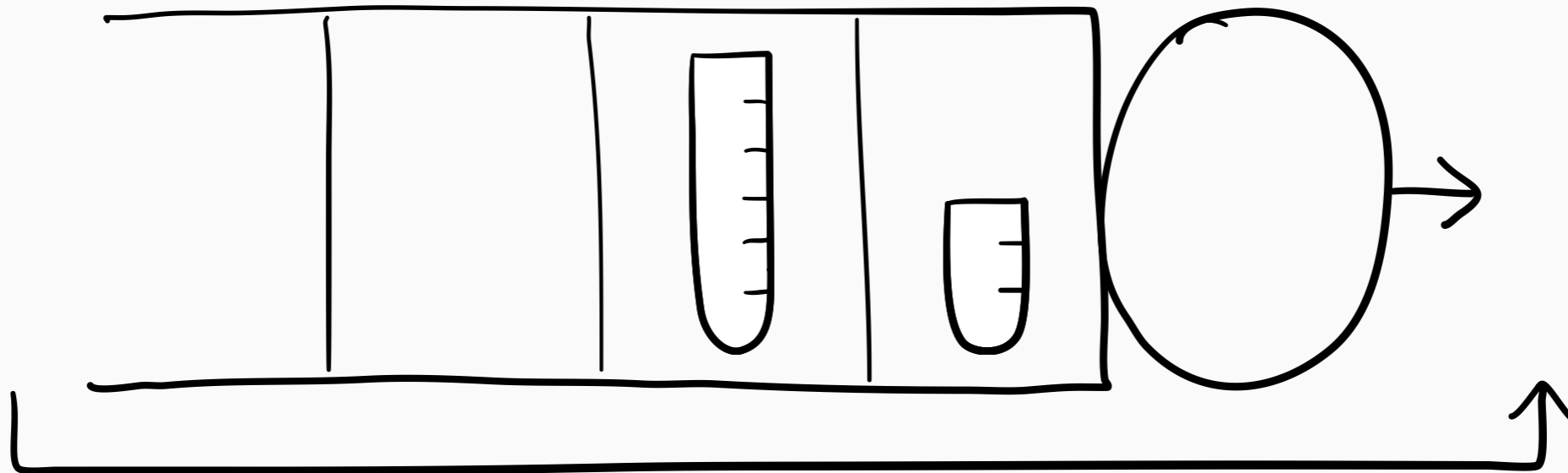
# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$



$T$ = response time

# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$

$T$ = response time

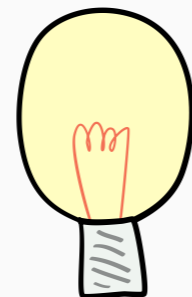# How should we schedule jobs to minimize delay?



$\mathbf{E}[T]$

Serve short jobs before long jobs

$T$ = response time

# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$

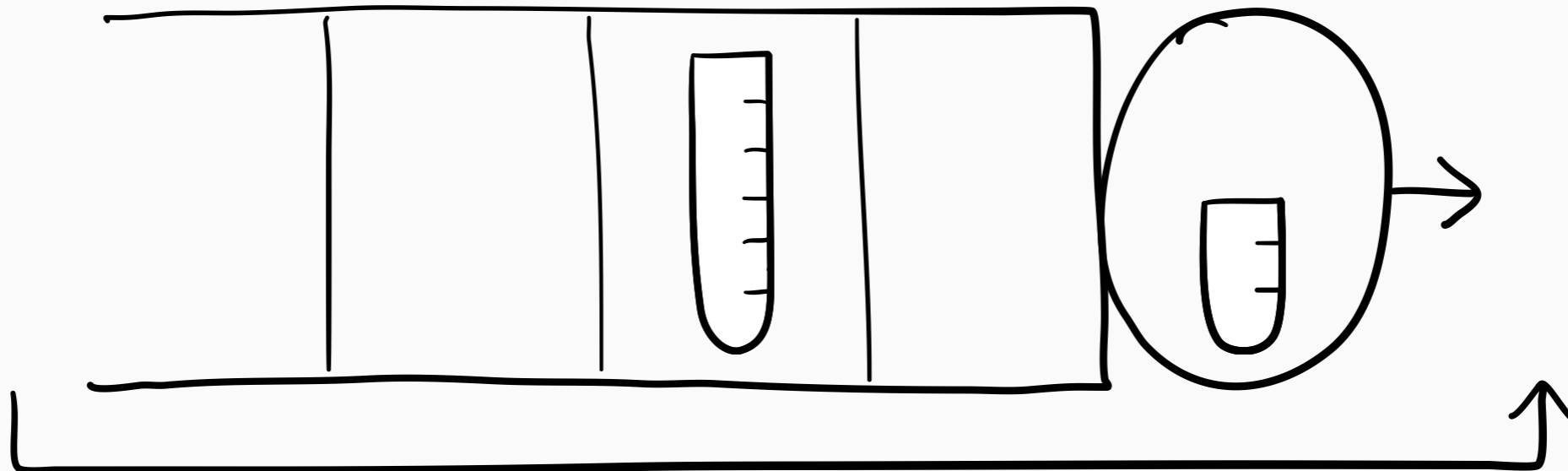Serve short jobs before long jobs

$T$ = response time

# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$

Serve short jobs before long jobs

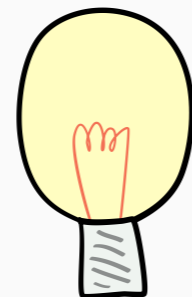$T$ = response time

# How should we schedule jobs to minimize delay?
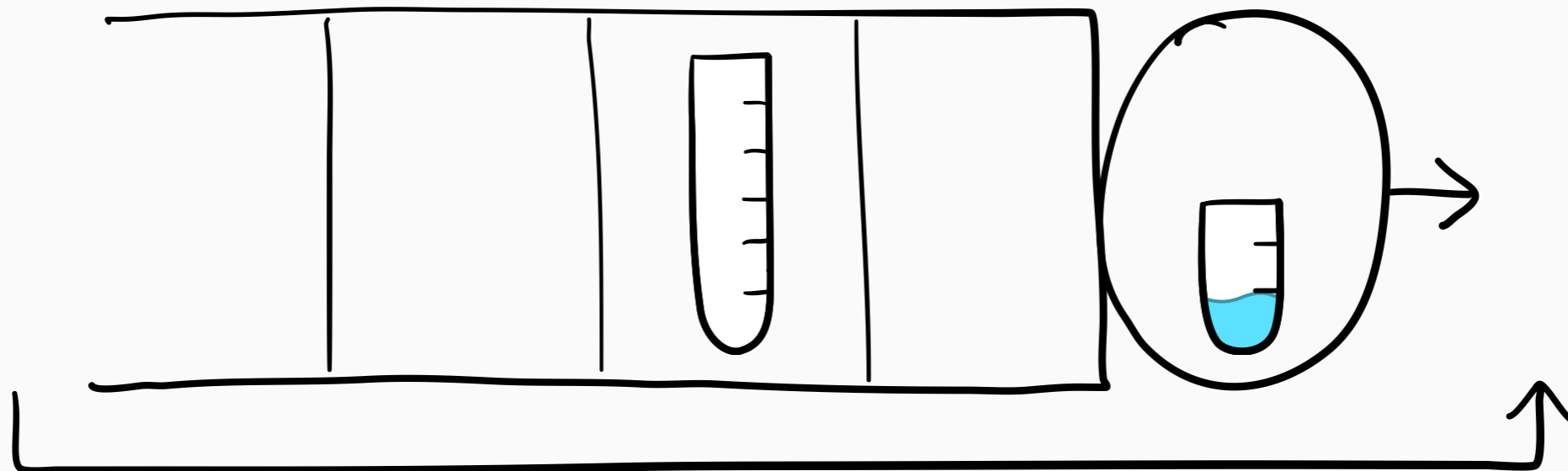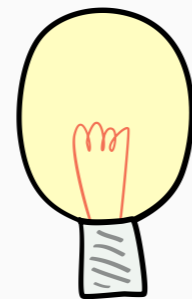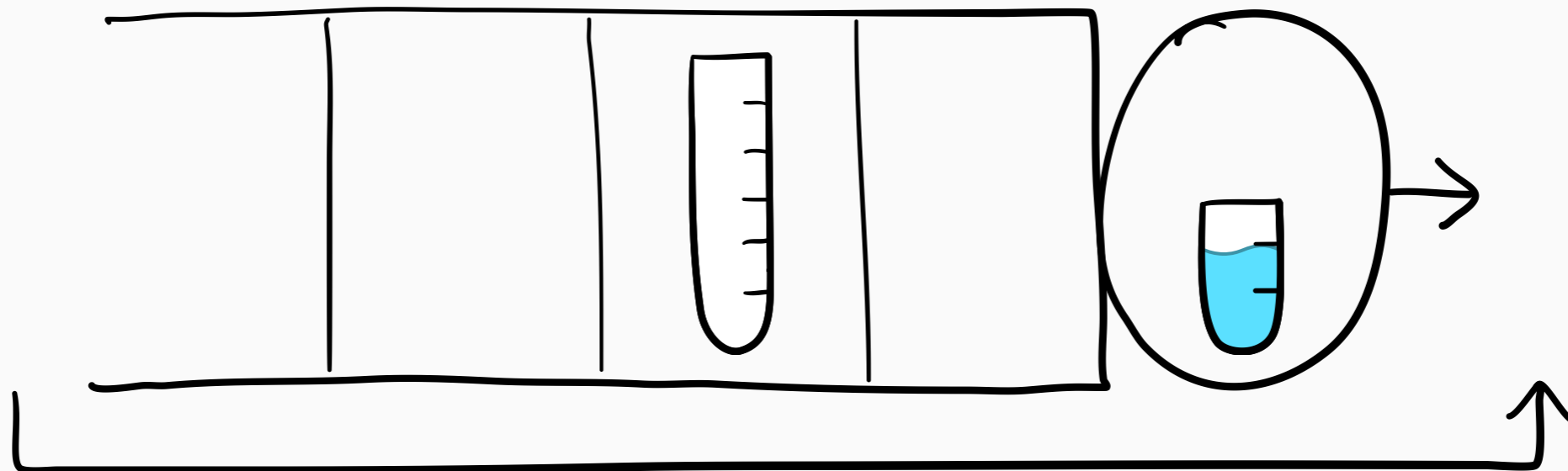
$\mathbf{E}[T]$

Serve short jobs before long jobs

$T$ = response time

# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$

Serve short jobs before long jobs

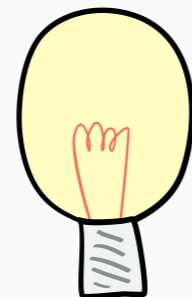$T$ = response time

# How should we schedule jobs to minimize delay?
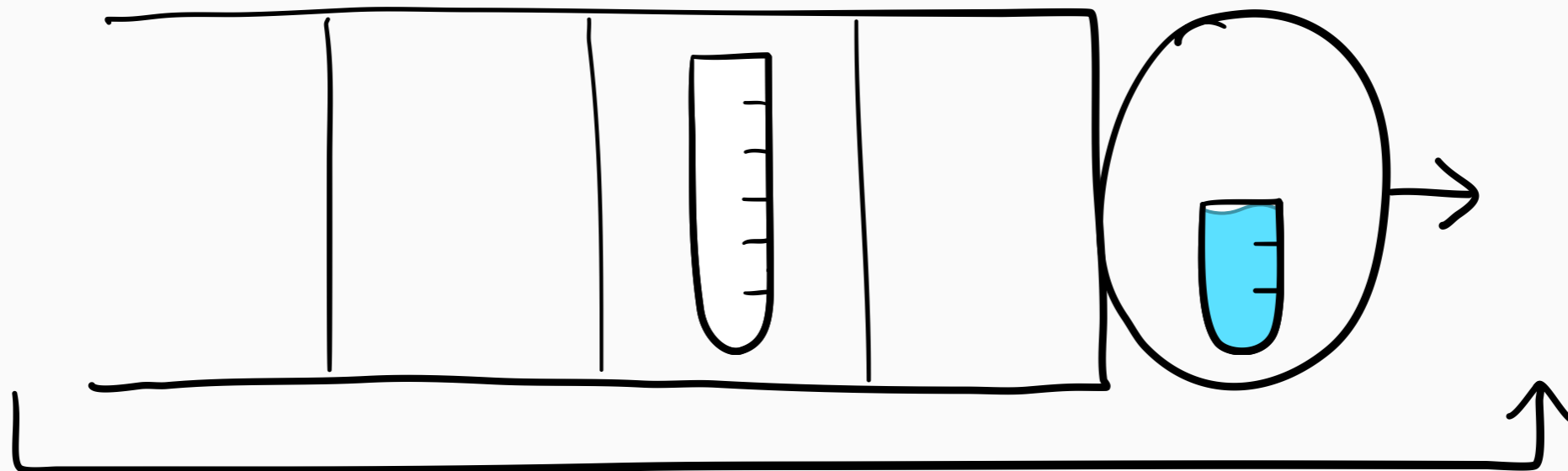
$\mathbf{E}[T]$

Serve short jobs before long jobs

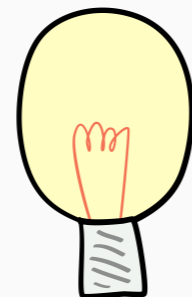$T$ = response time

# How should we schedule jobs to minimize delay?
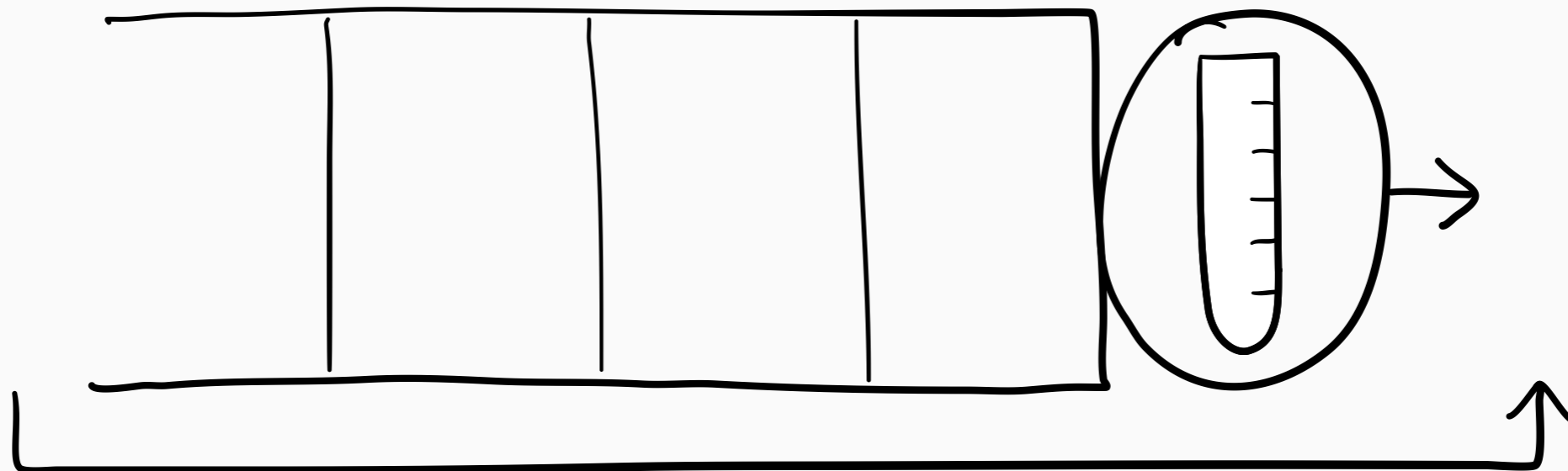
$\mathbf{E}[T]$

Serve short jobs before long jobs

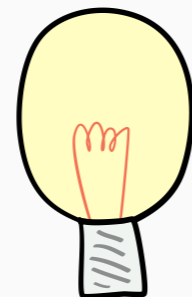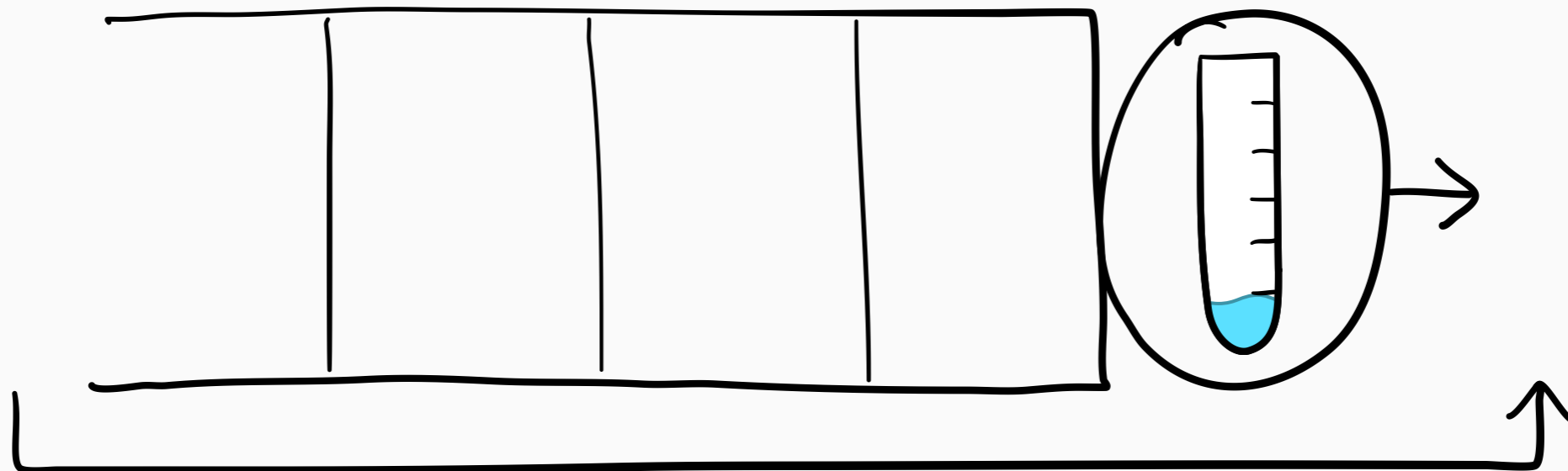$T$ = response time

# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$

Serve short jobs before long jobs

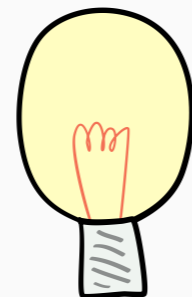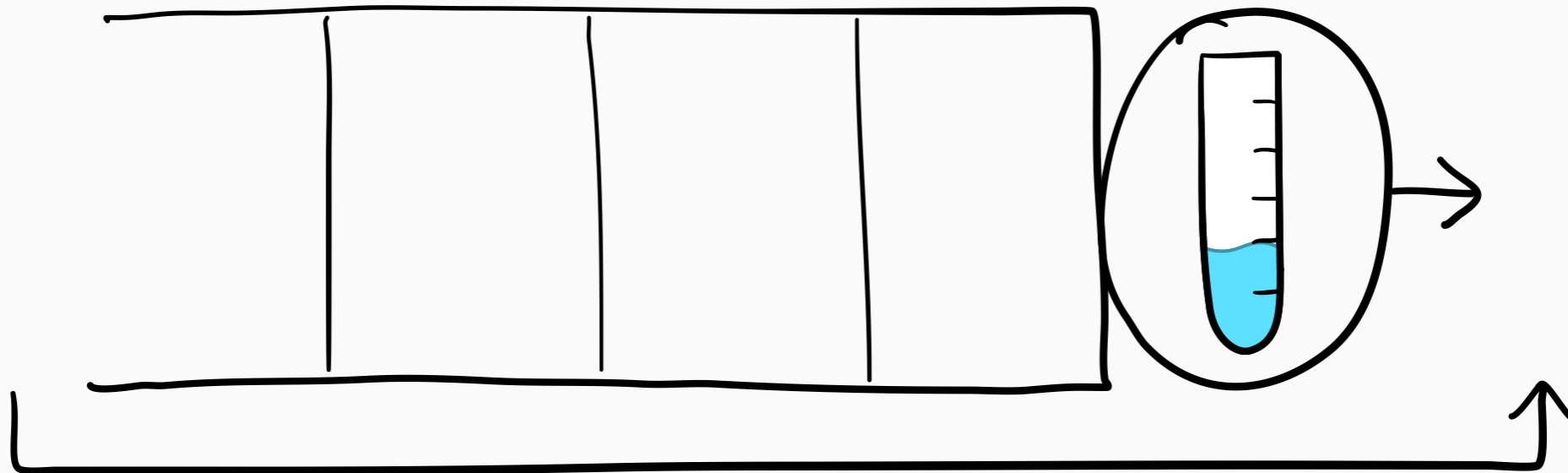$T$ = response time

# How should we schedule jobs to minimize delay?

$\mathbf{E}[T]$

Serve short jobs before long jobs

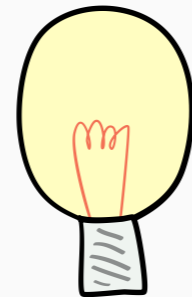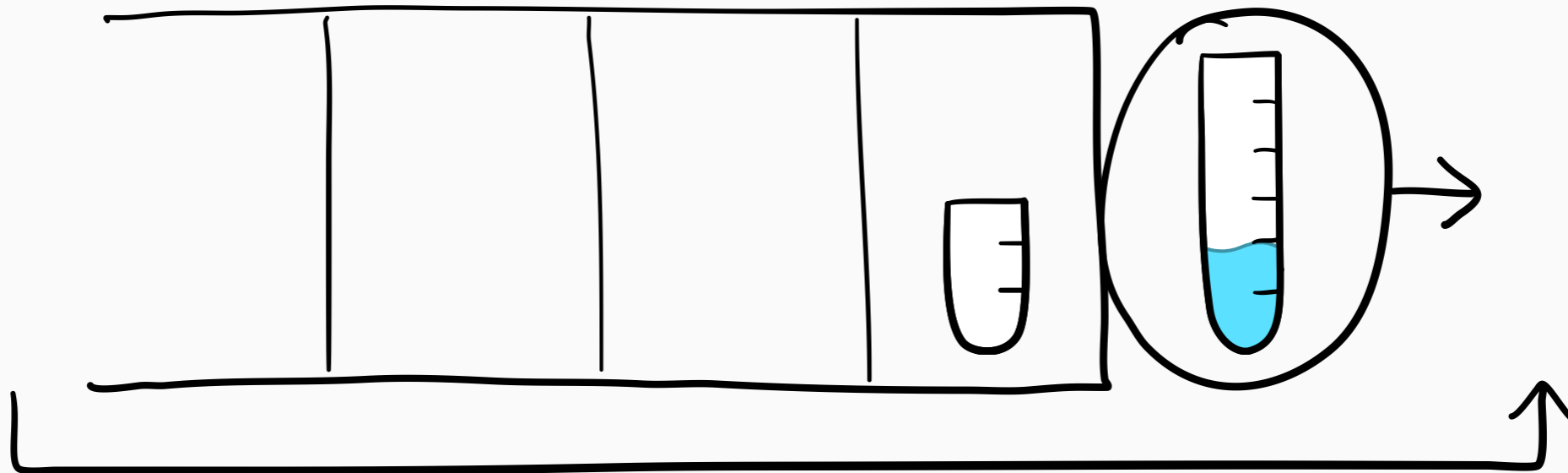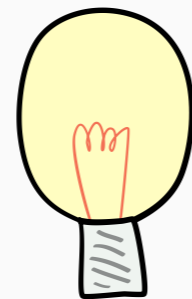$T$ = response time

# How should we schedule jobs to minimize delay?

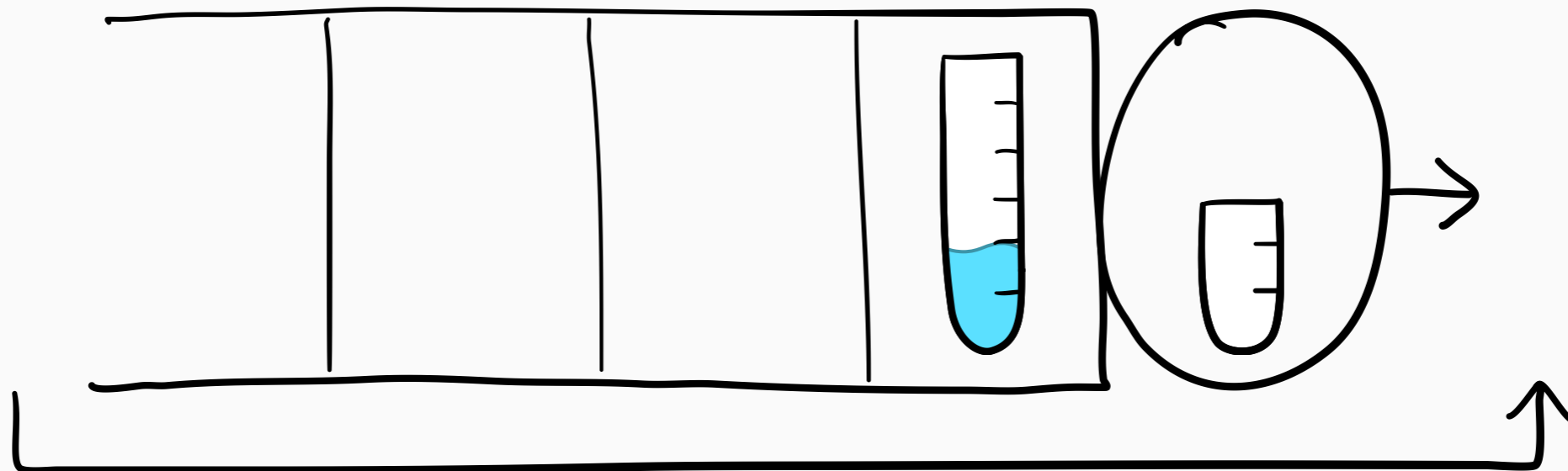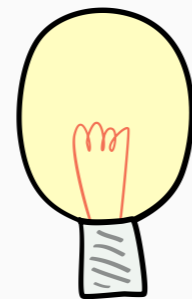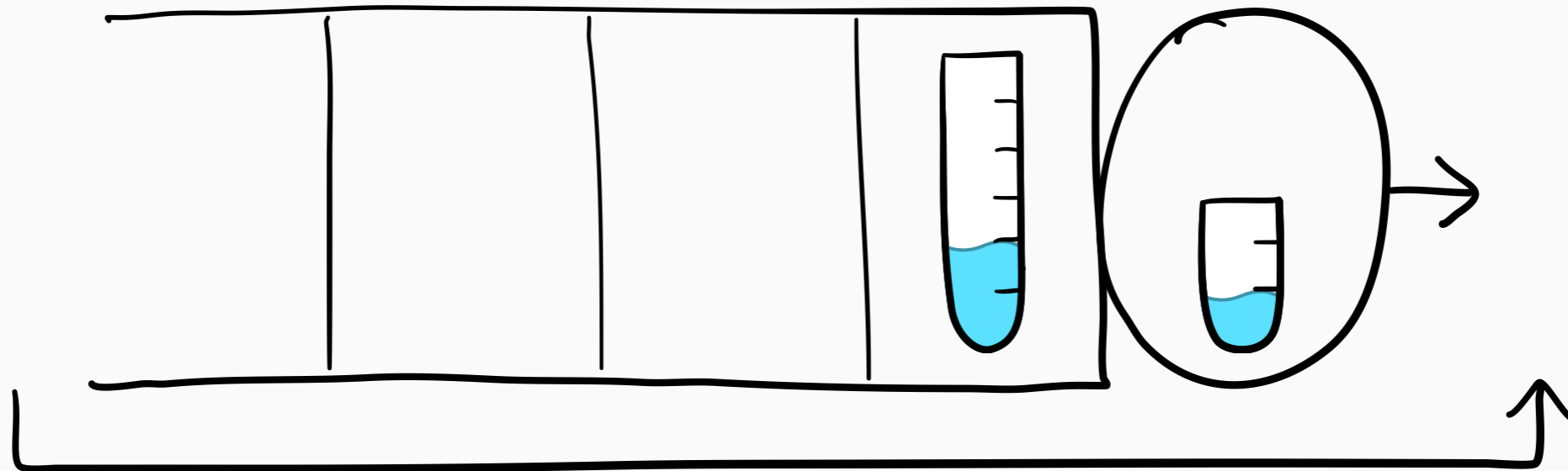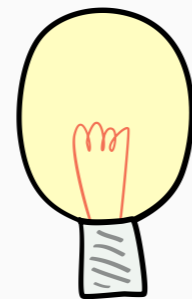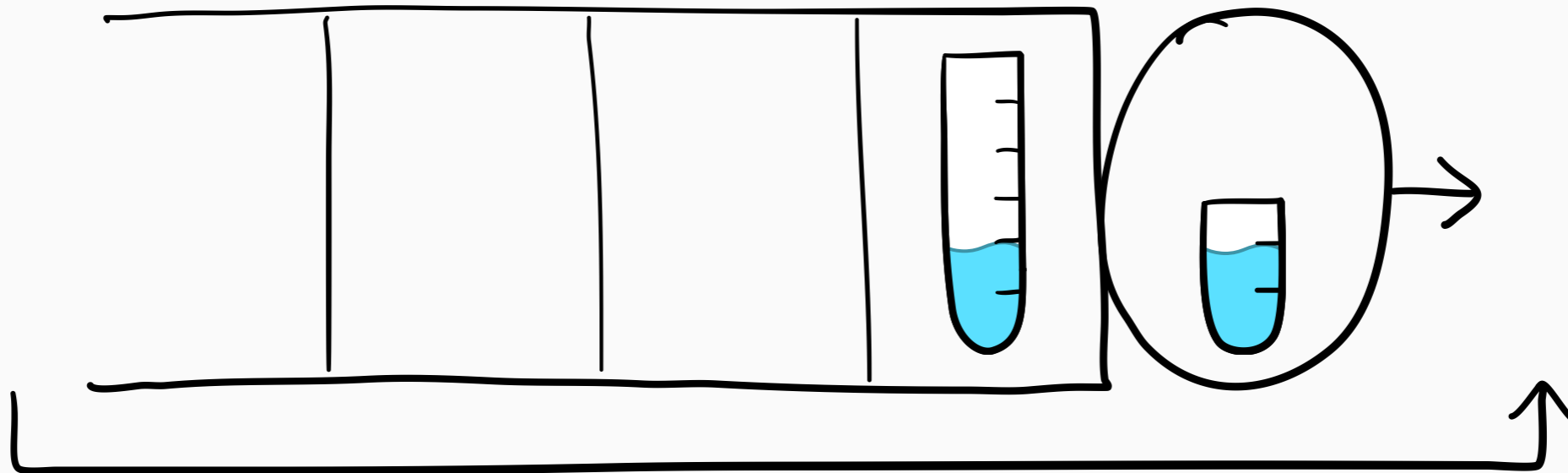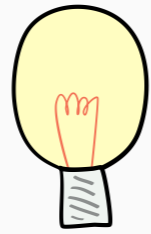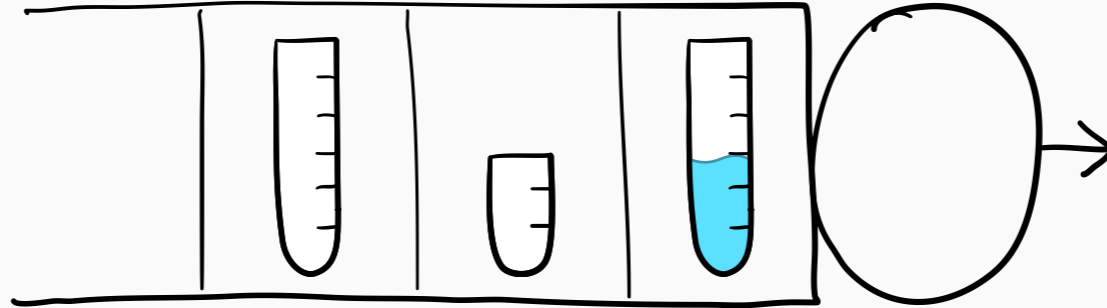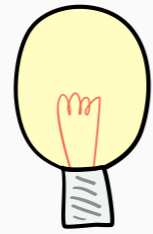$\mathbf{E}[T]$

Serve short jobs before long jobs

$T$ = response time

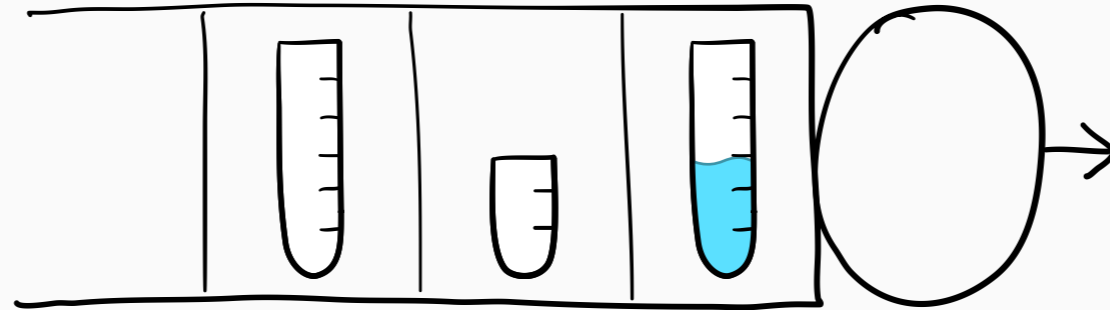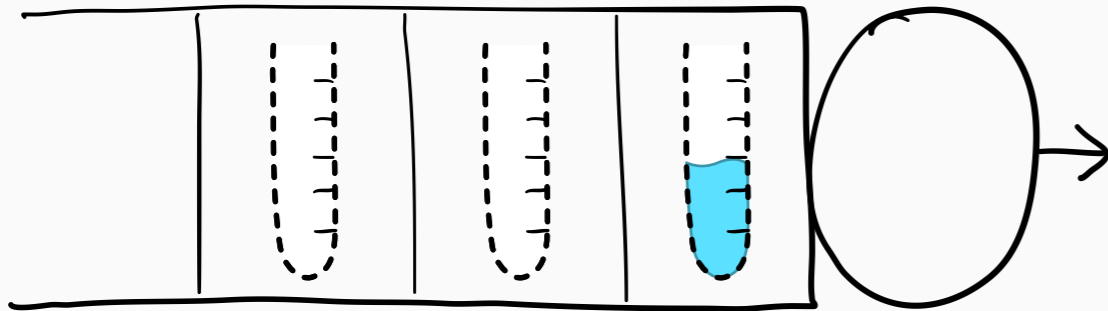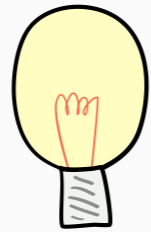short before long



3

short before long



**unknown** sizes

short before long

unknown sizes

multiple servers

$k$

short before long

**unknown** sizes

**multiple** servers

$k$

**This work: both** at once!

SRPT
(Schrage, 1968)

short before long

⚠ **unknown** sizes

⚠ **multiple** servers

$k$

⚠⚠ **This work: both at once!**

SRPT (Schrage, 1968)

short before long

SRPT-*k* (GSH, 2018)

⚠ **unknown** sizes

⚠ **multiple** servers

*k*

⚠⚠ **This work: both at once!**

SRPT
(Schrage, 1968)

short before long

**Gittins**
(several, 1970s)

**unknown** sizes

SRPT-*k*
(GSH, 2018)

**multiple** servers

*k*

**This work: both at once!**

3

# Main result

**Theorem:** **Gittins**-***k*** has "near-optimal" $\mathbf{E}[T]$ in the M/G/***k*** with **unknown** job sizes



$T$ = response time

# Main result

**Theorem:** **Gittins-*k*** has "near-optimal" $\mathbf{E}[T]$ in the M/G/*k* with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$



$T$ = response time

# Main result

**Theorem: Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
**A:** Use **Gittins**



$T$ = response time

# Main result

**Theorem:** **Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
**A:** Use **Gittins**



$T$ = response time

# Main result

**Theorem:** **Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
**A:** Use **Gittins**



$T$ = response time

# Main result

**Theorem: Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
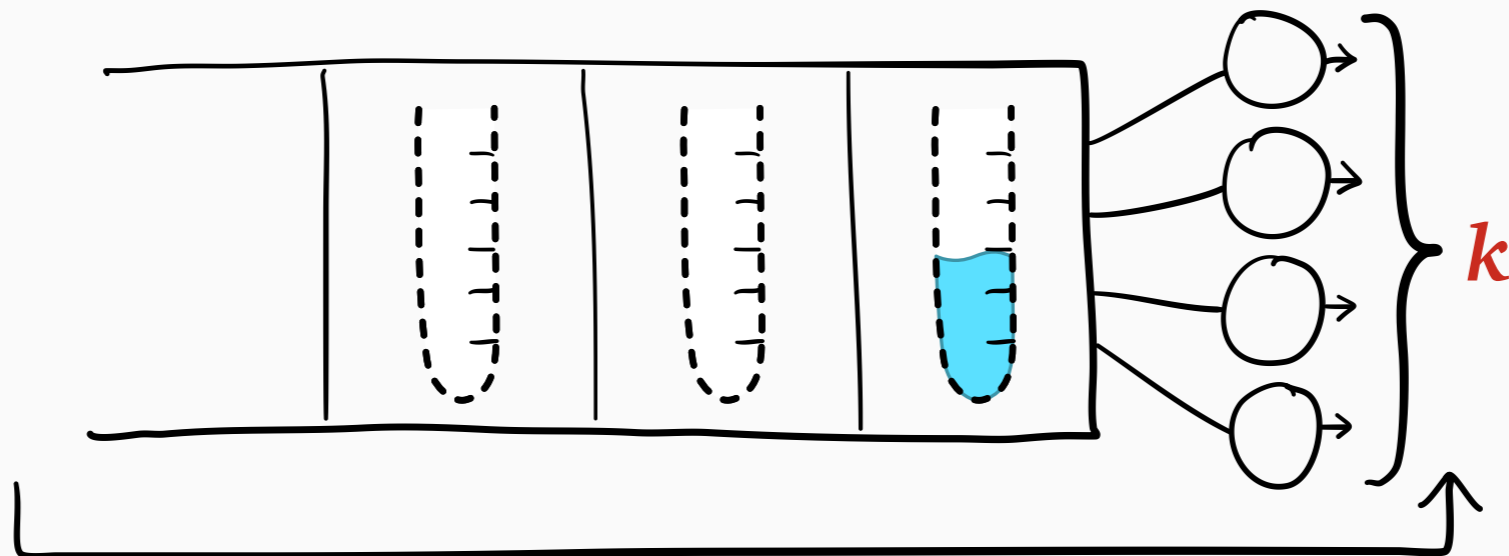**A:** Use **Gittins**



$T$ = response time

4

# Main result

**Theorem:** **Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes
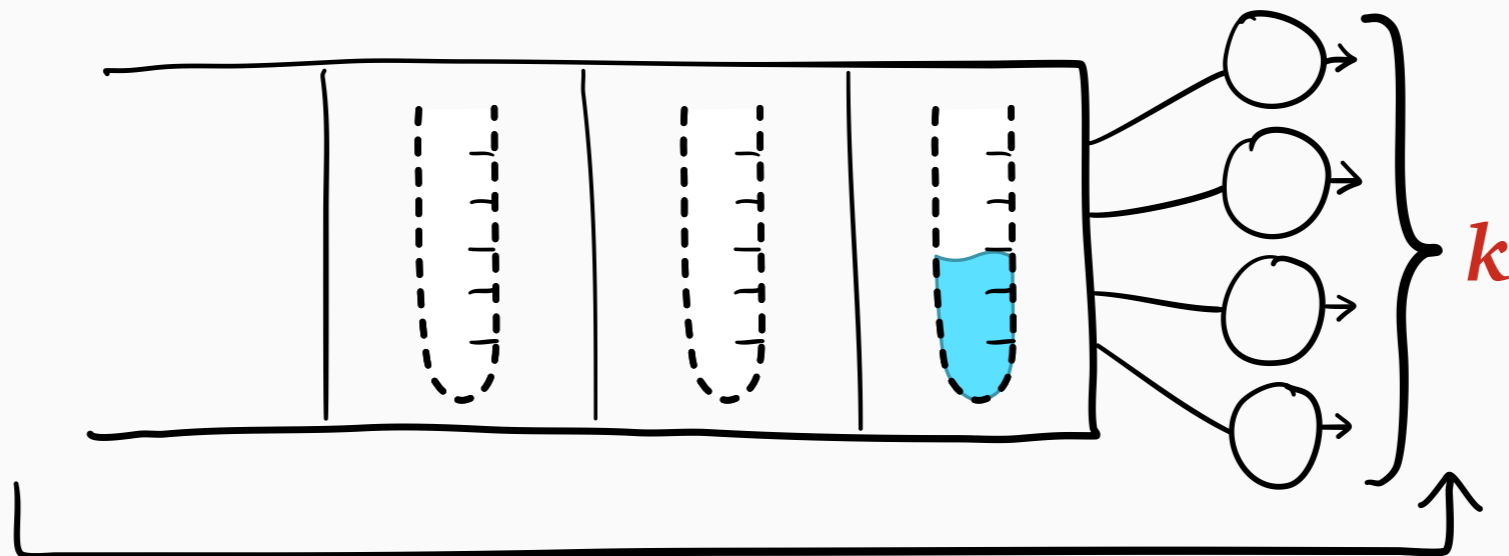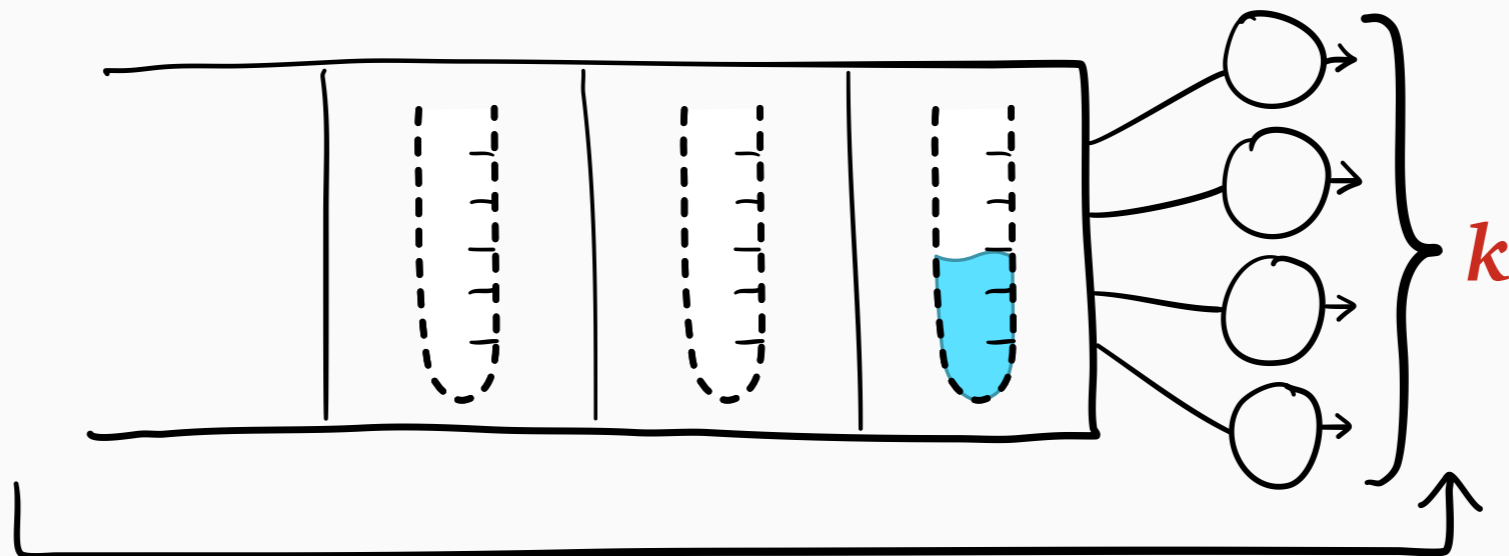
$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
**A:** Use **Gittins**



$T$ = response time

# Main result

**Theorem:** **Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$
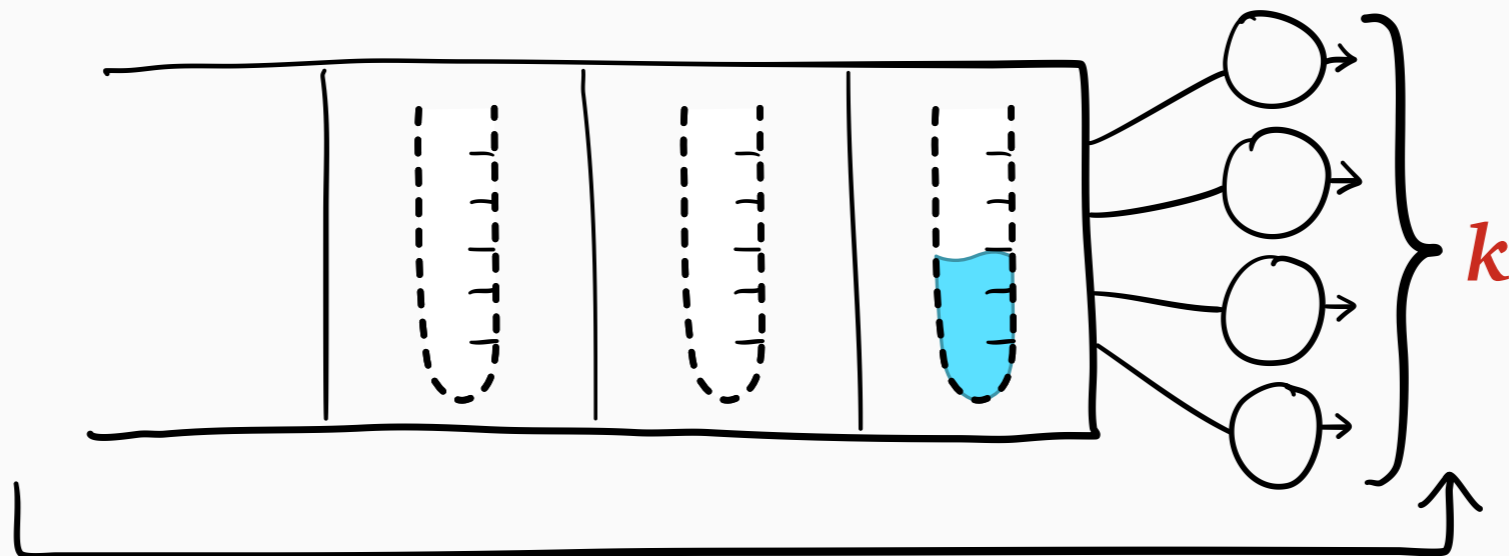in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
**A:** Use **Gittins**

# Main result

**Theorem:** **Gittins**-$k$ has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
**A:** Use **Gittins**

**Q:** How do we analyze $\mathbf{E}[T]$?

# Main result

**Theorem:** **Gittins-$k$** has "near-optimal" $\mathbf{E}[T]$ in the M/G/$k$ with **unknown** job sizes

$$\mathbf{E}[T_{\text{Gittins-}k}] \leq \mathbf{E}[T_{\text{Opt-}k}] + \text{"small"}$$

**Q:** How to schedule with $k$ servers?
 **A:** Use **Gittins**

**Q:** How do we analyze $\mathbf{E}[T]$?
 **Q:** Why is our approach significant?

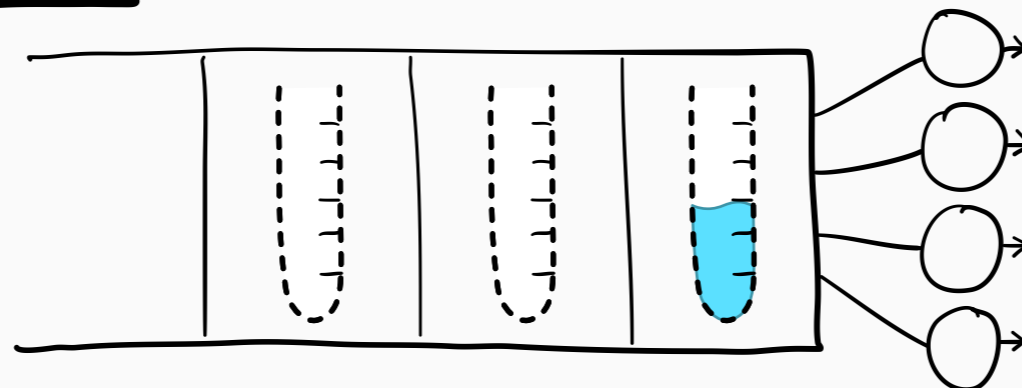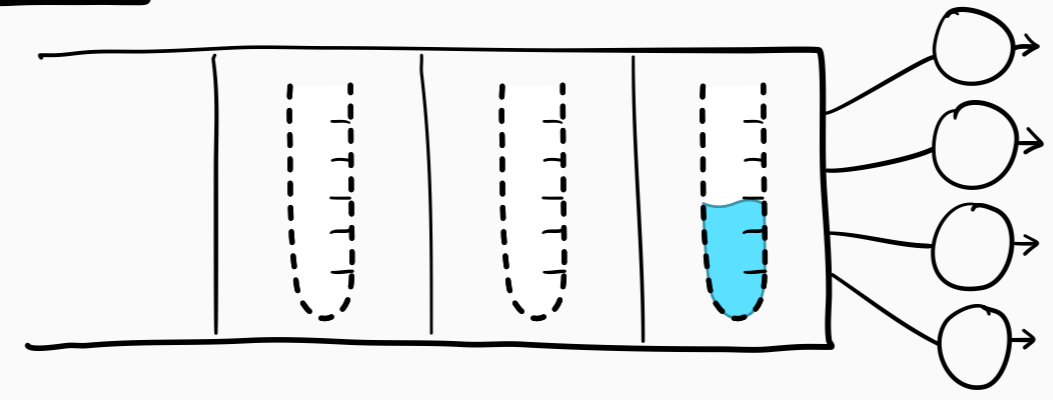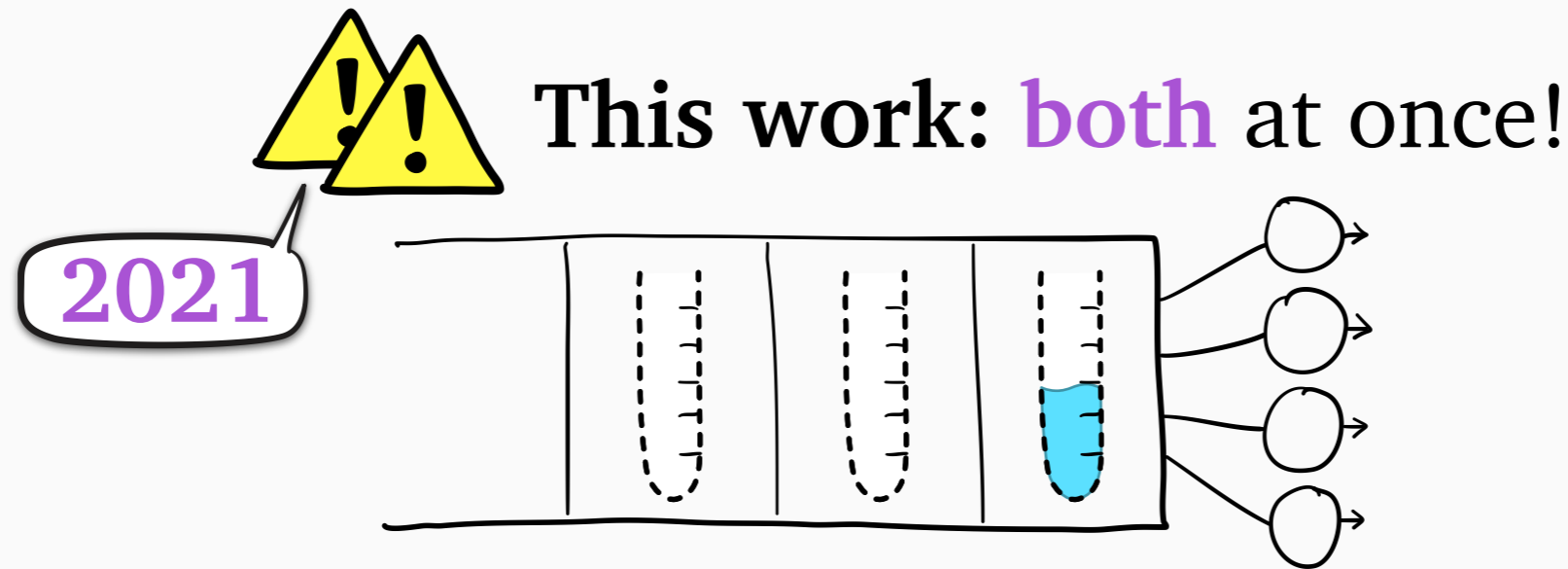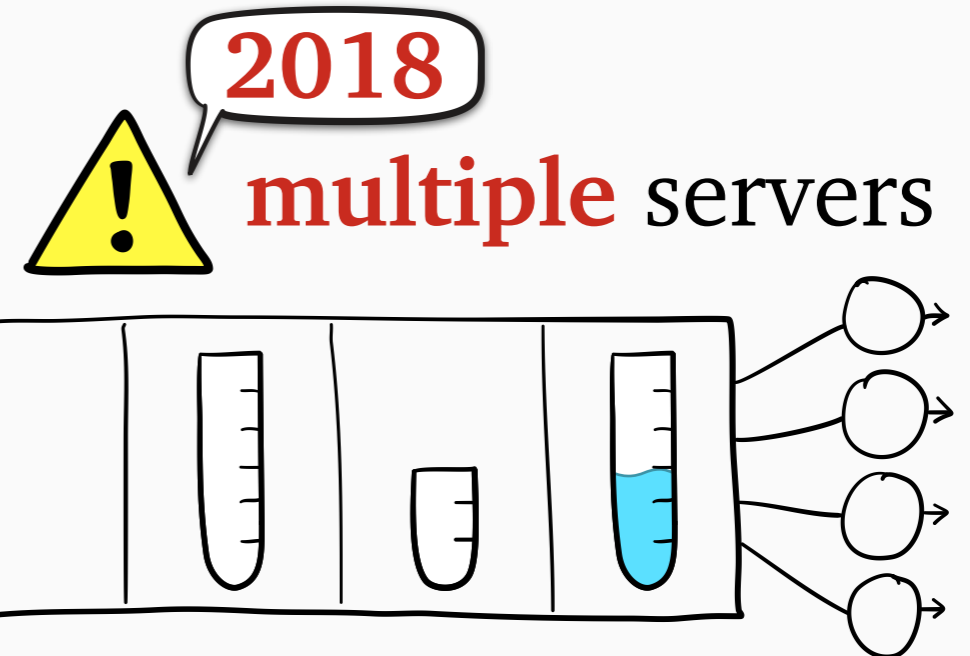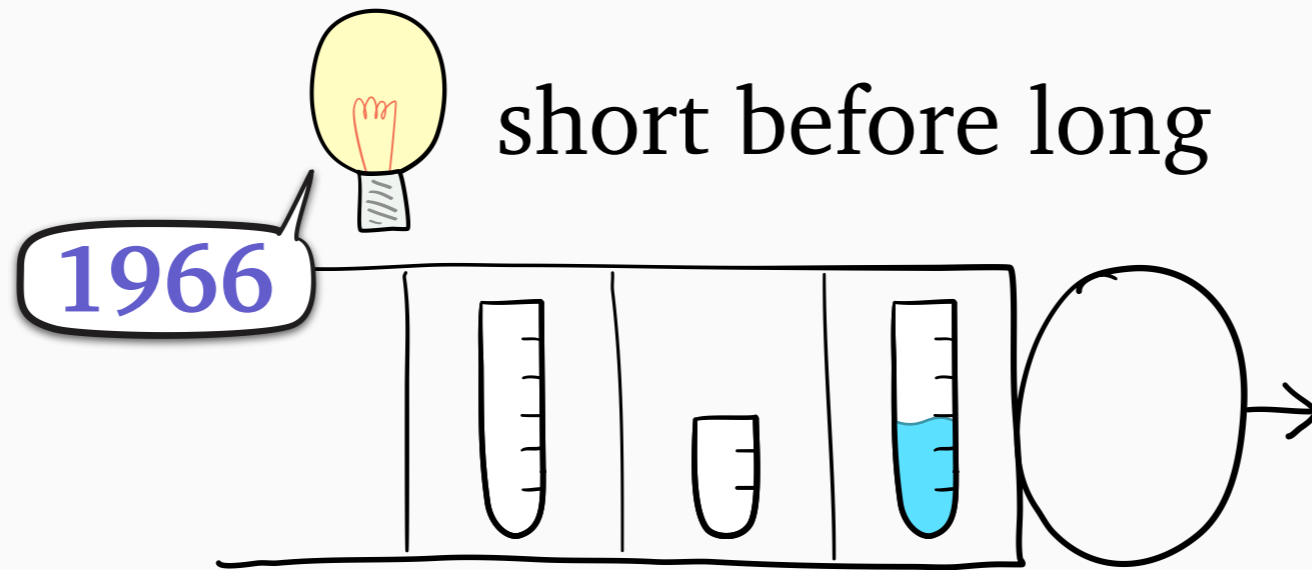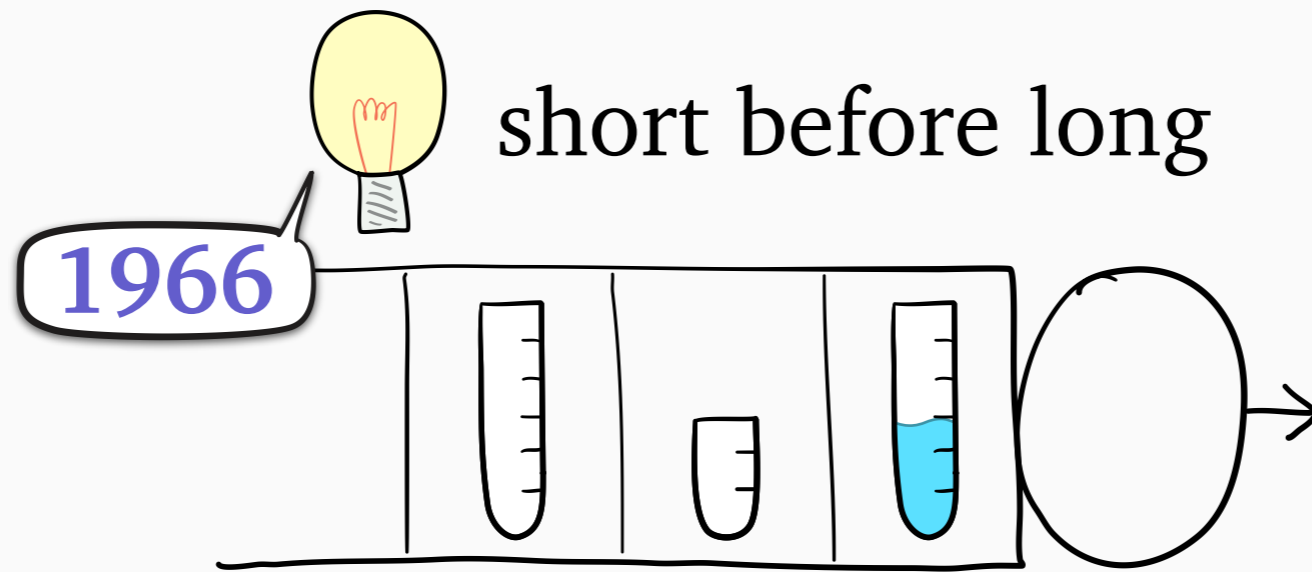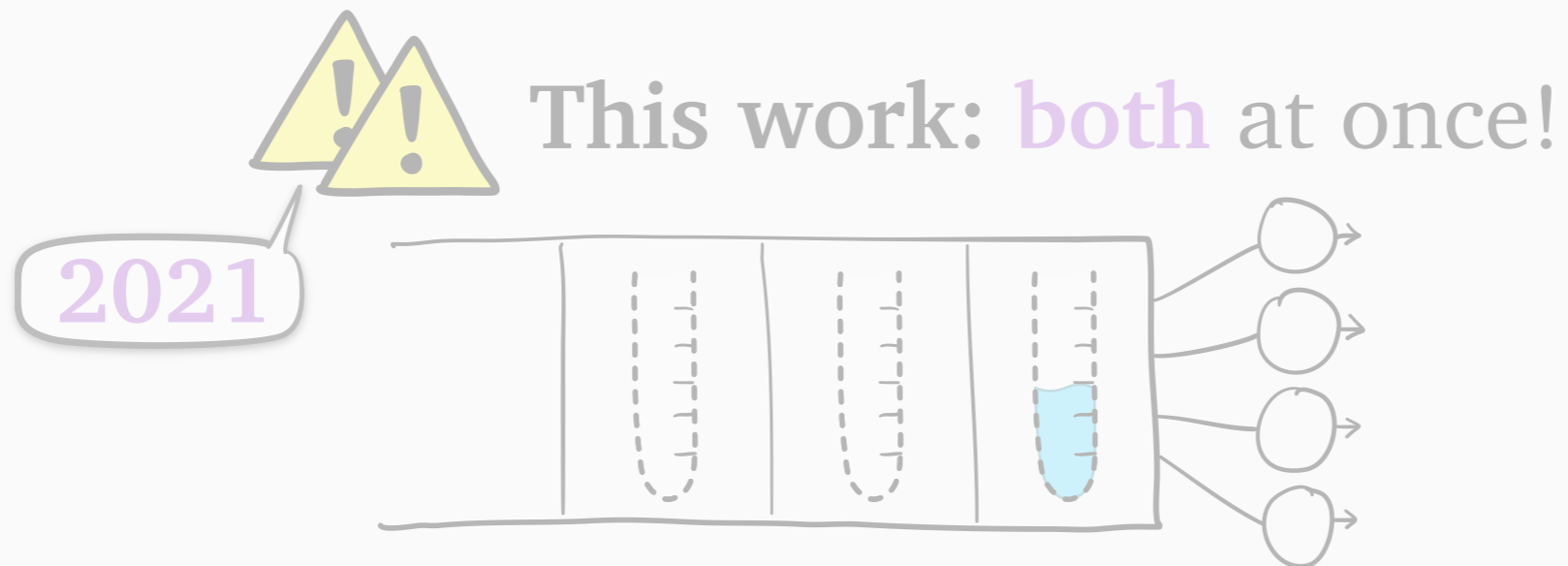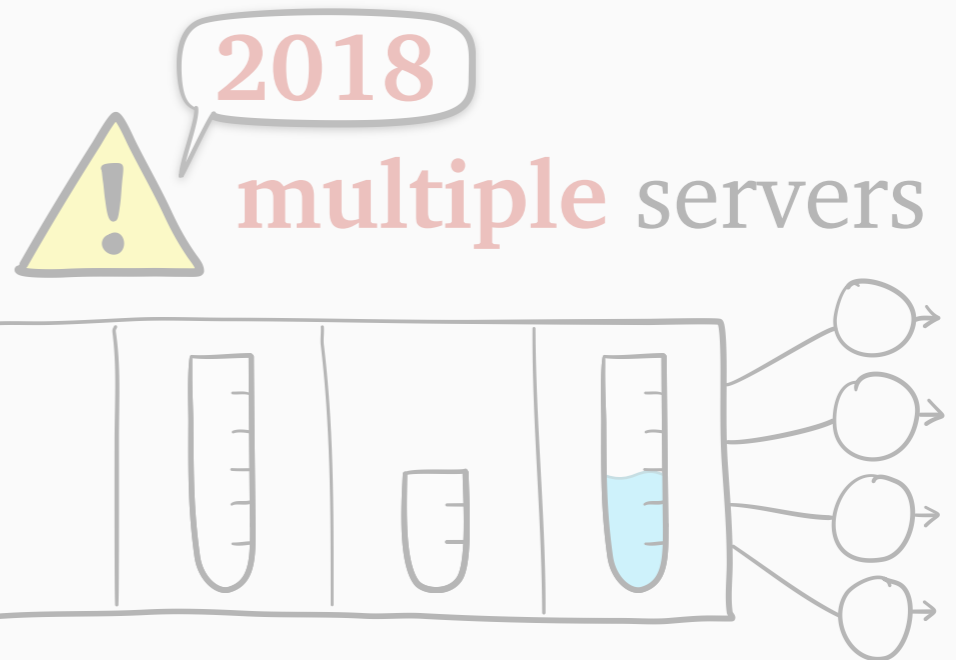short before long

SRPT
(Schrage, 1968)

Gittins
(several, 1970s)

unknown sizes

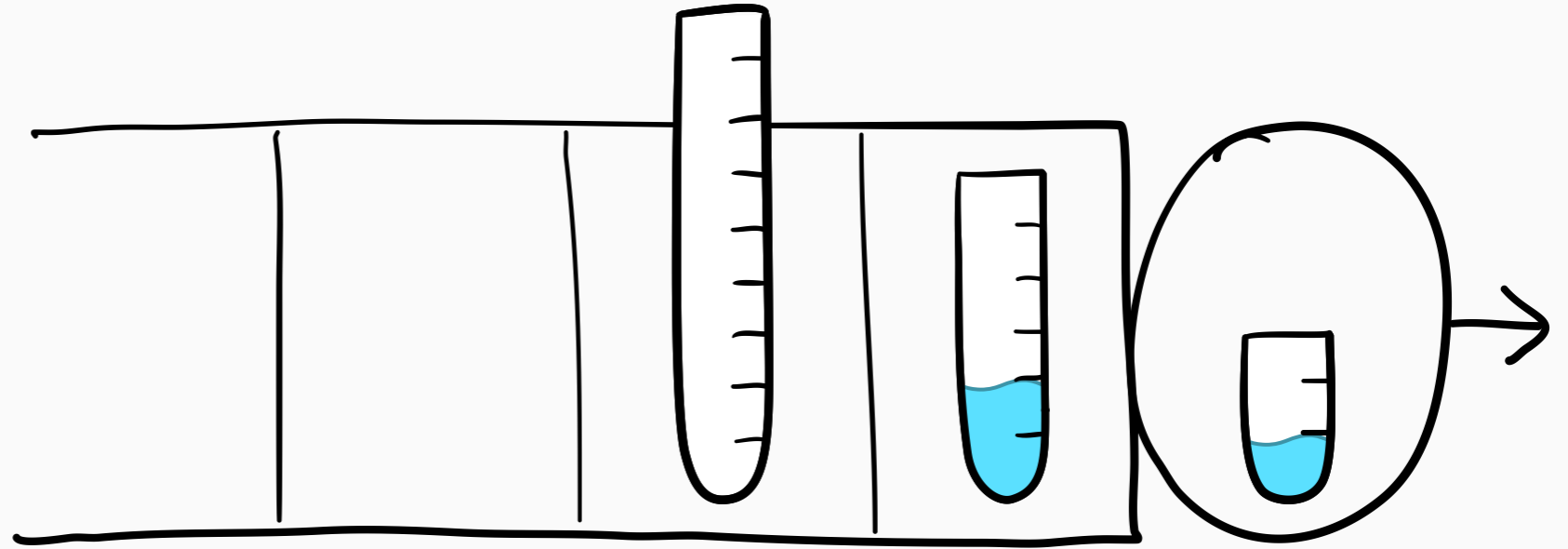SRPT-$k$
(GSH, 2018)

multiple servers

This work: both at once!

5

# $\mathbf{E}[T]$ of SRPT

# E[$T$] of SRPT



**tagged job**

# $\mathrm{E}[T]$ of SRPT



**tagged job**

random system state

# $\mathbf{E}[T]$ of SRPT



**tagged job**   | = **rank**   random system state

remaining size, lower is better

# E[T] of SRPT



**tagged job**    ▯ = **rank**    random system state

remaining size, lower is better

# $\mathbf{E}[T]$ of SRPT



**tagged job**   = **rank**   random system state

remaining size, lower is better

**Key quantity:**
$W(r)$ = "$r$-work" = work relevant to job of **rank** $r$

# $\mathbf{E}[T]$ of SRPT



**tagged job**    = **rank**    random system state

remaining size, lower is better

**Key quantity:**
$W(r)$ = "$r$-work" = work relevant to job of **rank** $r$

6

# E[$T$] of SRPT



**tagged job**    ▯ = **rank**    rank-system state

remaining size,
lower is better

**Key quantity:**
   $W(r) = $ "$r$-work" = work relevant to job of **rank** $r$

# Two-step $\mathrm{E}[T]$ analysis

# Two-step $\mathbf{E}[T]$ analysis

**Step 1:** compute $\mathbf{E}[W(r)]$

# Two-step $\mathbf{E}[T]$ analysis

**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$

**Step 2:** $\mathbf{E}[W(\textcolor{purple}{r})]$ to $\mathbf{E}[T]$

# Two-step $\mathrm{E}[T]$ analysis



**Step 1:** compute $\mathrm{E}[W(\textcolor{violet}{r})]$

**Step 2:** $\mathrm{E}[W(\textcolor{violet}{r})]$ to $\mathrm{E}[T]$

**Tagged job** method

# Two-step $\mathbf{E}[T]$ analysis

**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$

standard queueing

**Step 2:** $\mathbf{E}[W(\textcolor{purple}{r})]$ to $\mathbf{E}[T]$

**Tagged job** method

Tagged job methods

# $\mathbf{E}[T]$ with **multiple** servers

**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$

**Step 2:** $\mathbf{E}[W(\textcolor{purple}{r})]$ to $\mathbf{E}[T]$

# $\mathbf{E}[T]$ with **multiple** servers

**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$

SRPT-$\textcolor{red}{k}$

**Step 2:** $\mathbf{E}[W(\textcolor{purple}{r})]$ to $\mathbf{E}[T]$

# $\mathbf{E}[T]$ with **multiple** servers

**Step 1:** compute $\mathbf{E}[W(r)]$

SRPT-$k$

! intractable

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

# $\mathbf{E}[T]$ with **multiple** servers

**Step 1:** compute $\mathbf{E}[W(r)]$

SRPT-$k$

SRPT-1

⚠ intractable

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

# $\mathbf{E}[T]$ with **multiple** servers

# $\mathbf{E}[T]$ with **multiple** servers

**Step 1:** compute $\mathbf{E}[W(\mathbf{\color{purple}r})]$

SRPT-$\mathbf{\color{red}k}$ — speed $1/\mathbf{\color{red}k}$ ⟷ SRPT-1 — speed 1

⚠ intractable

**Worst-case** gap:

$$W_{\text{SRPT-}\mathbf{\color{red}k}}(\mathbf{\color{purple}r}) \leq W_{\text{SRPT-1}}(\mathbf{\color{purple}r}) + \mathbf{\color{red}k}\mathbf{\color{purple}r}$$

**Step 2:** $\mathbf{E}[W(\mathbf{\color{purple}r})]$ to $\mathbf{E}[T]$

# $\mathbf{E}[T]$ with **multiple** servers

**Step 1:** compute $\mathbf{E}[W(r)]$

SRPT-$k$    speed $1/k$      SRPT-1    speed 1

⚠ intractable

**Worst-case** gap:

$$W_{\text{SRPT-}k}(r) \leq W_{\text{SRPT-1}}(r) + kr$$
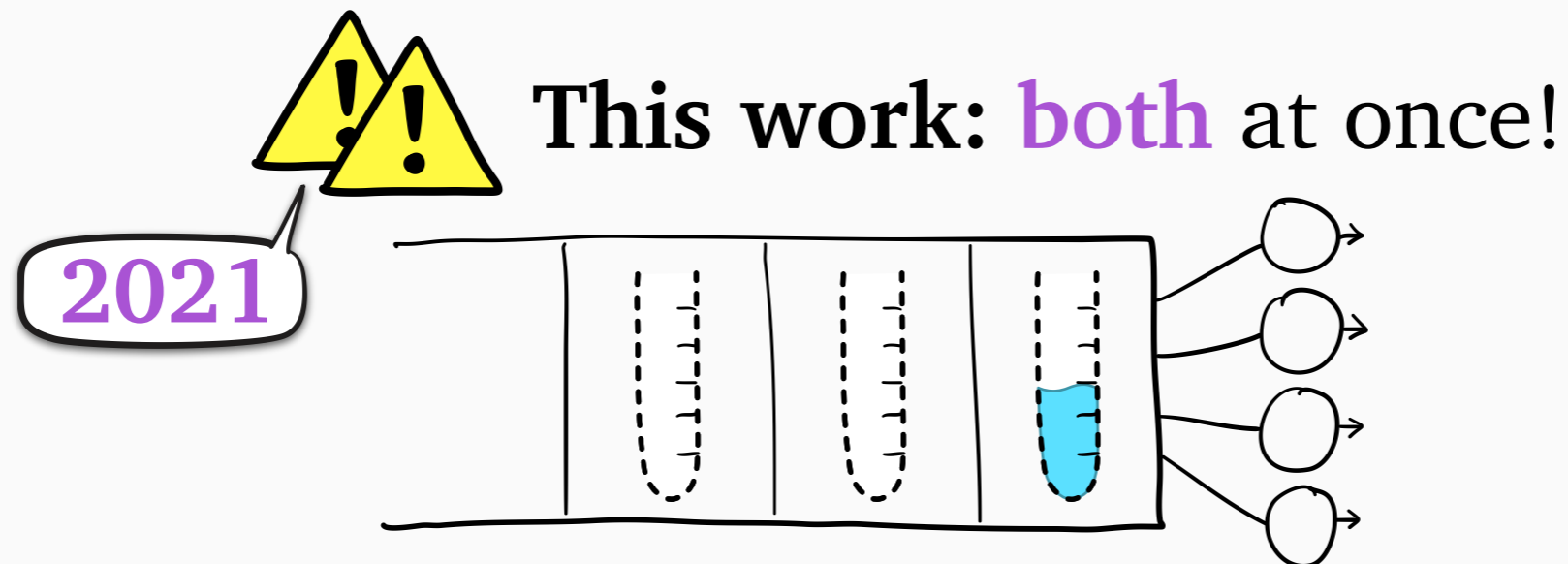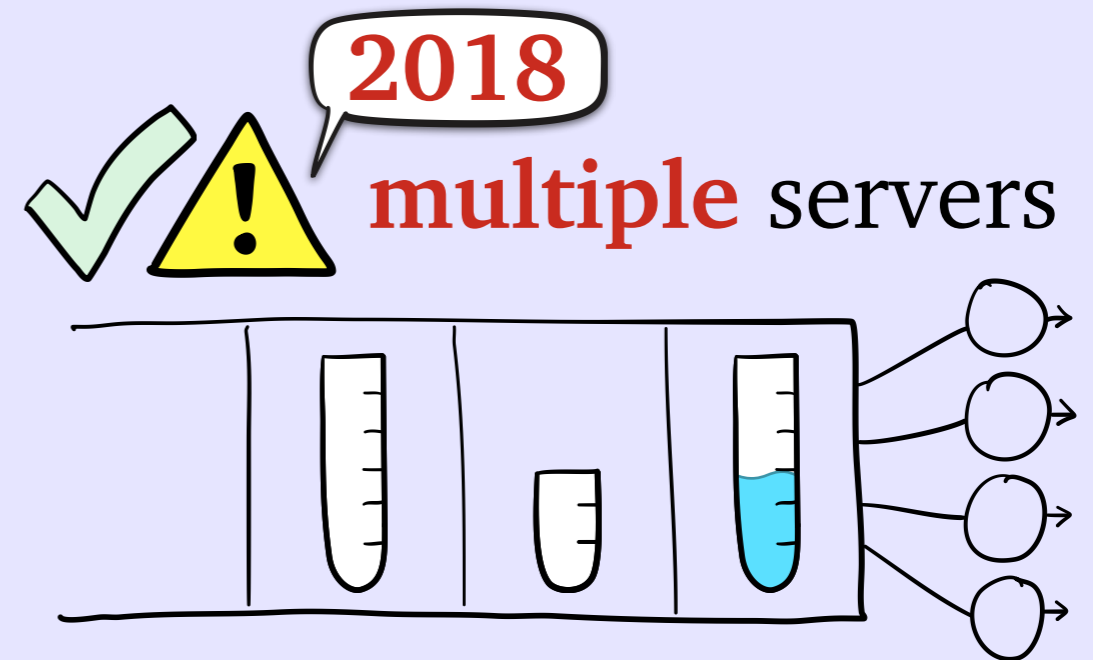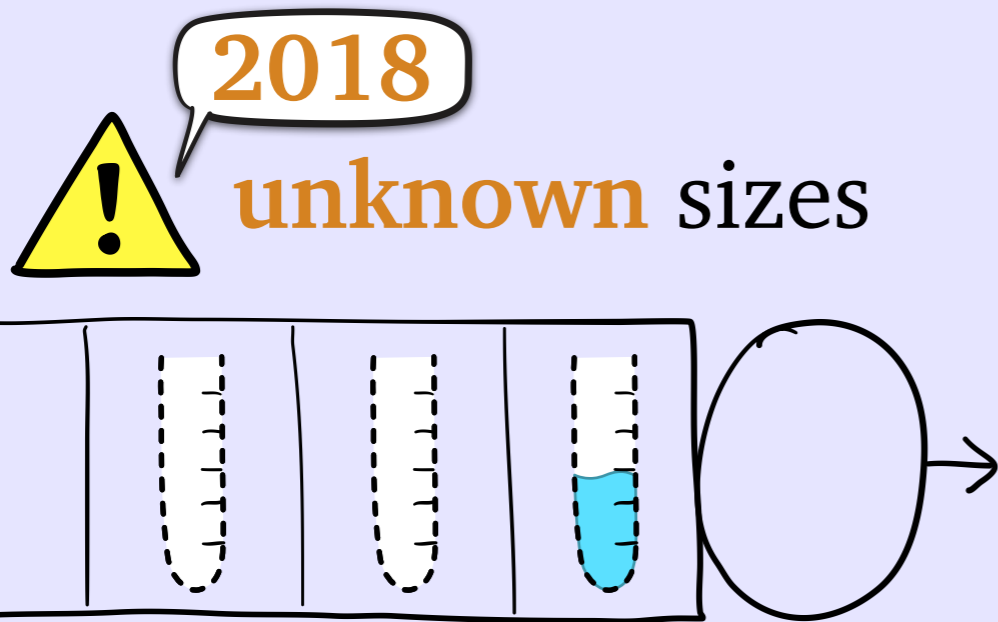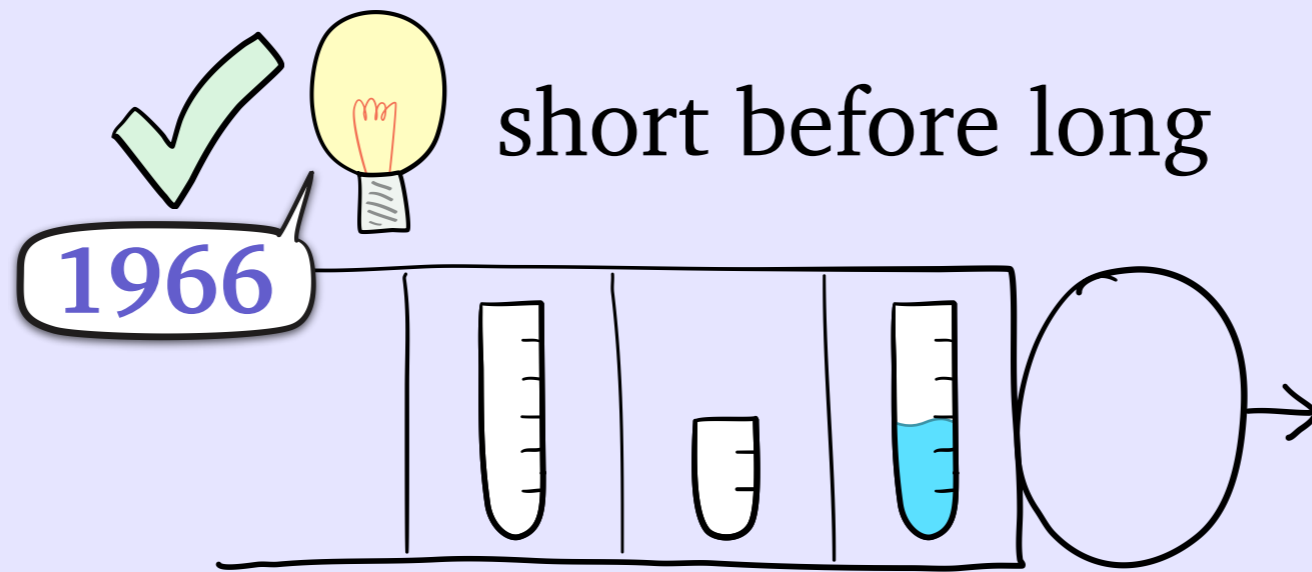
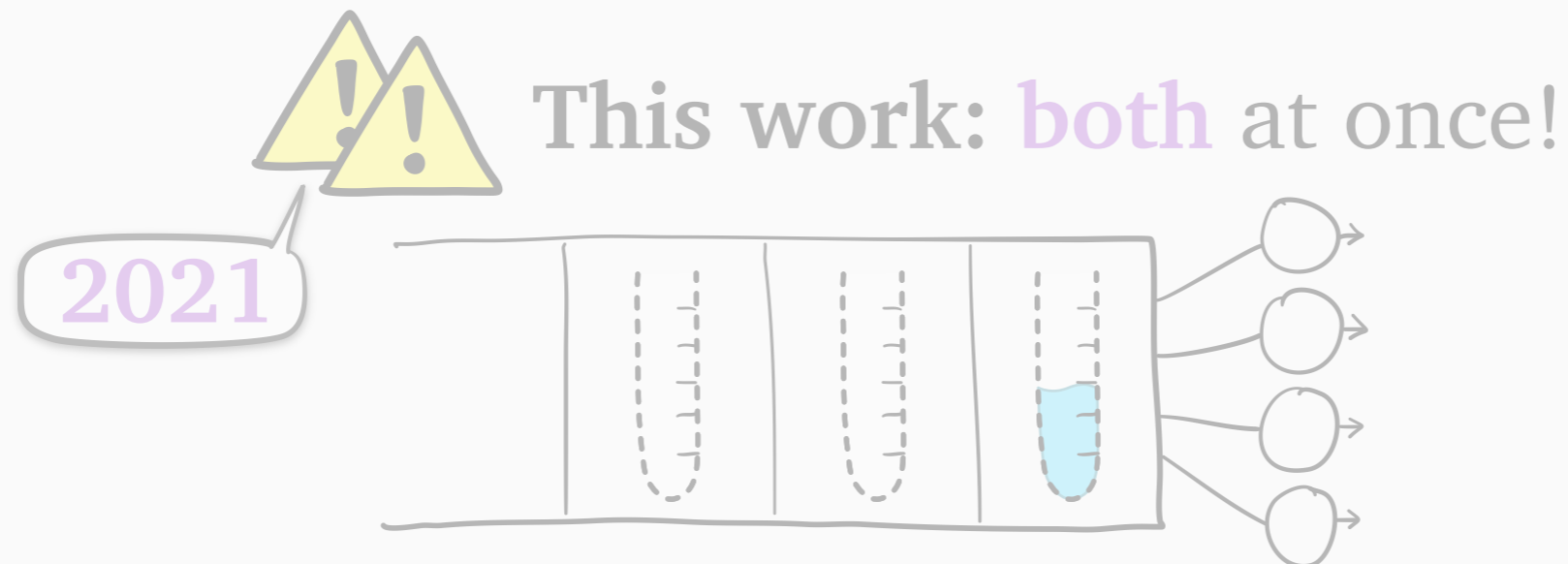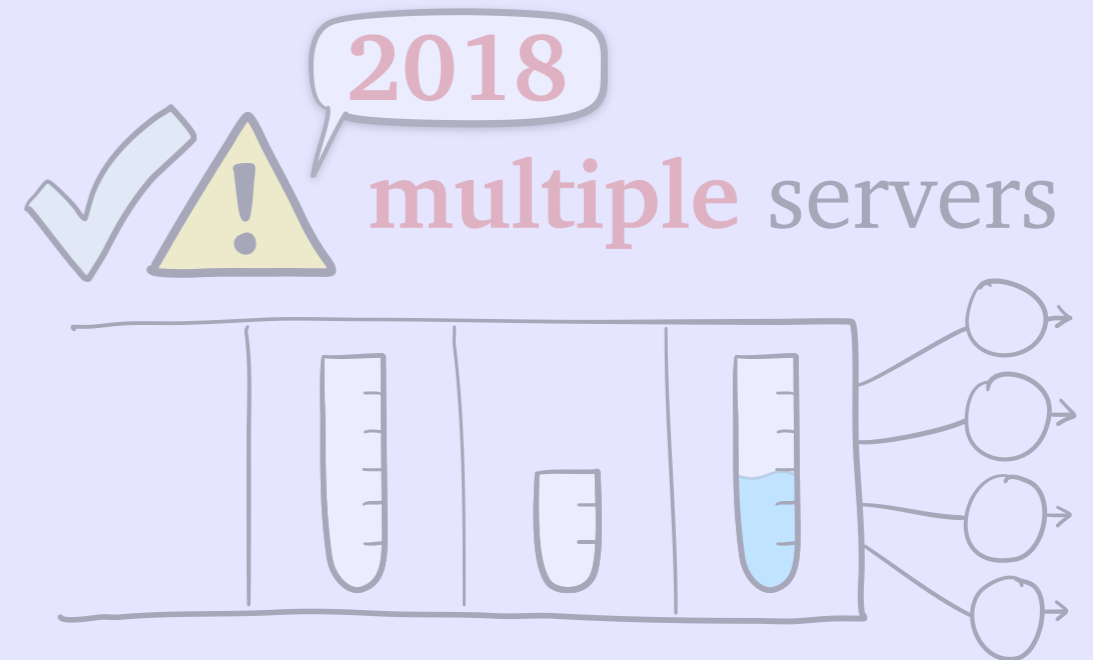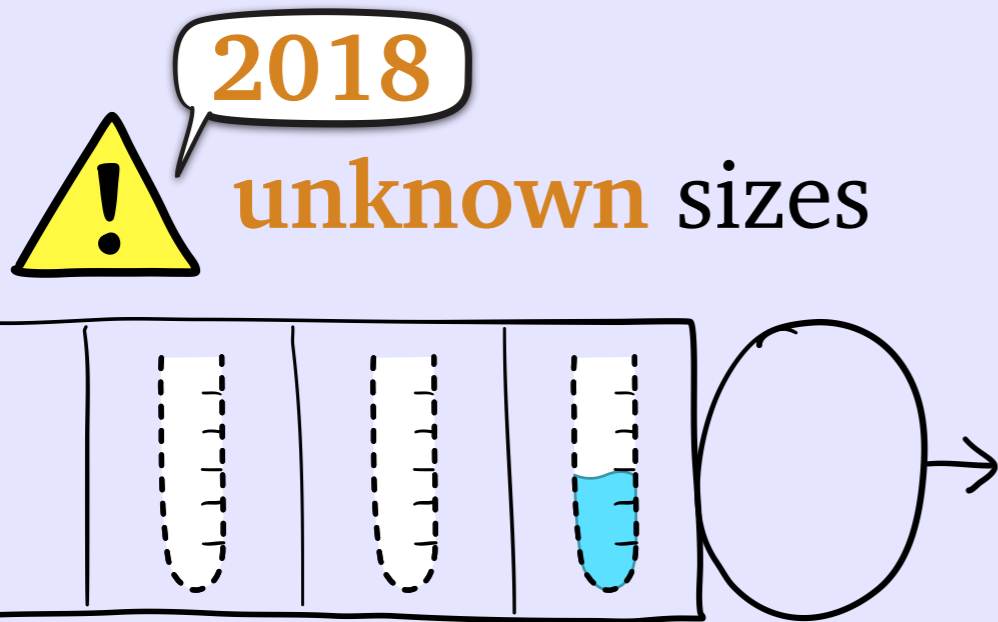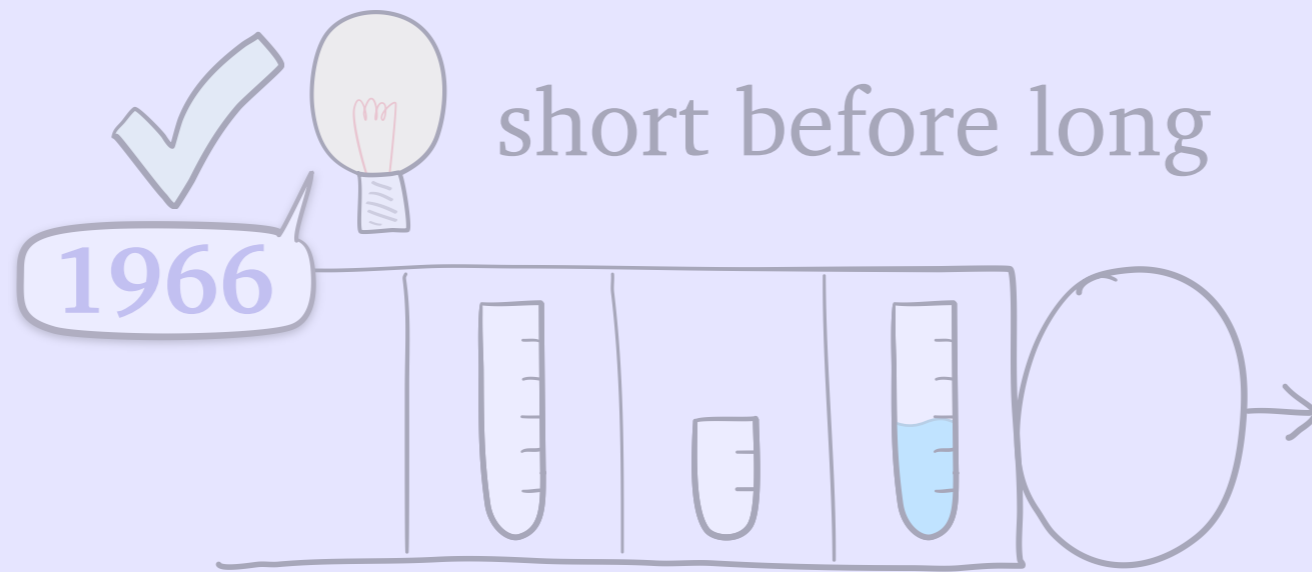**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

tagged job
+
worst-case
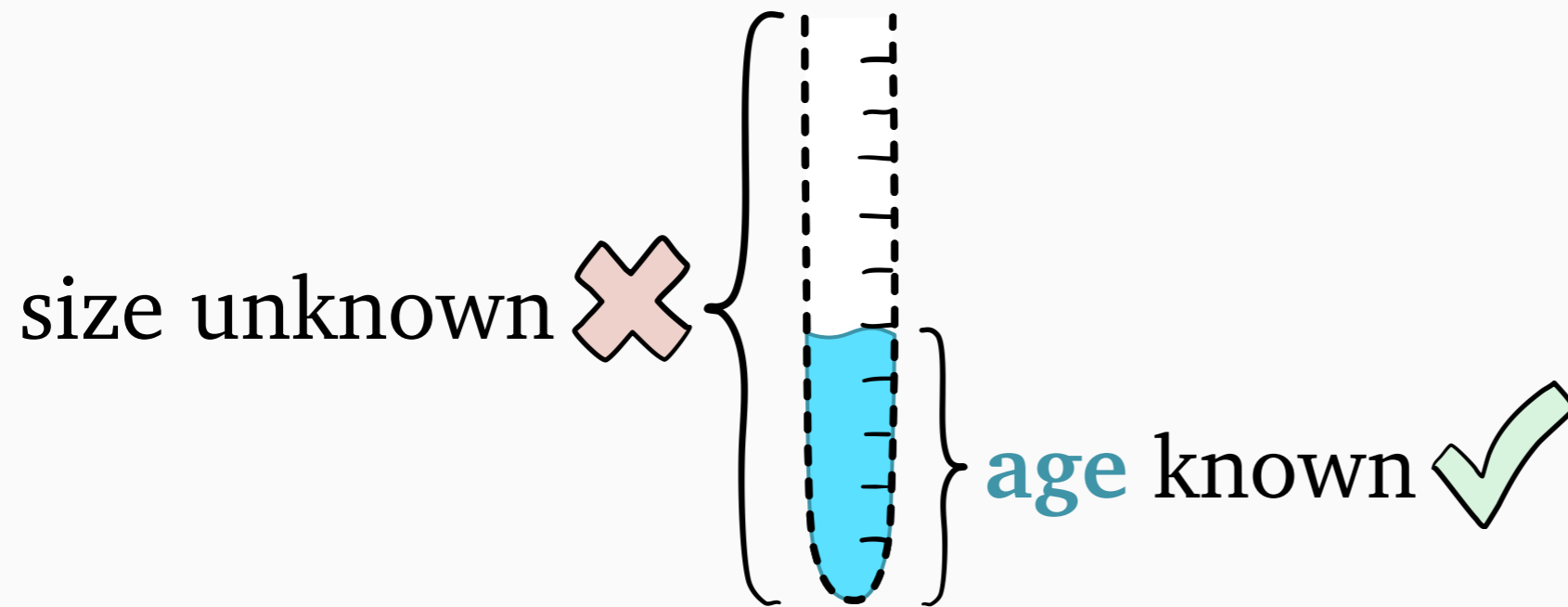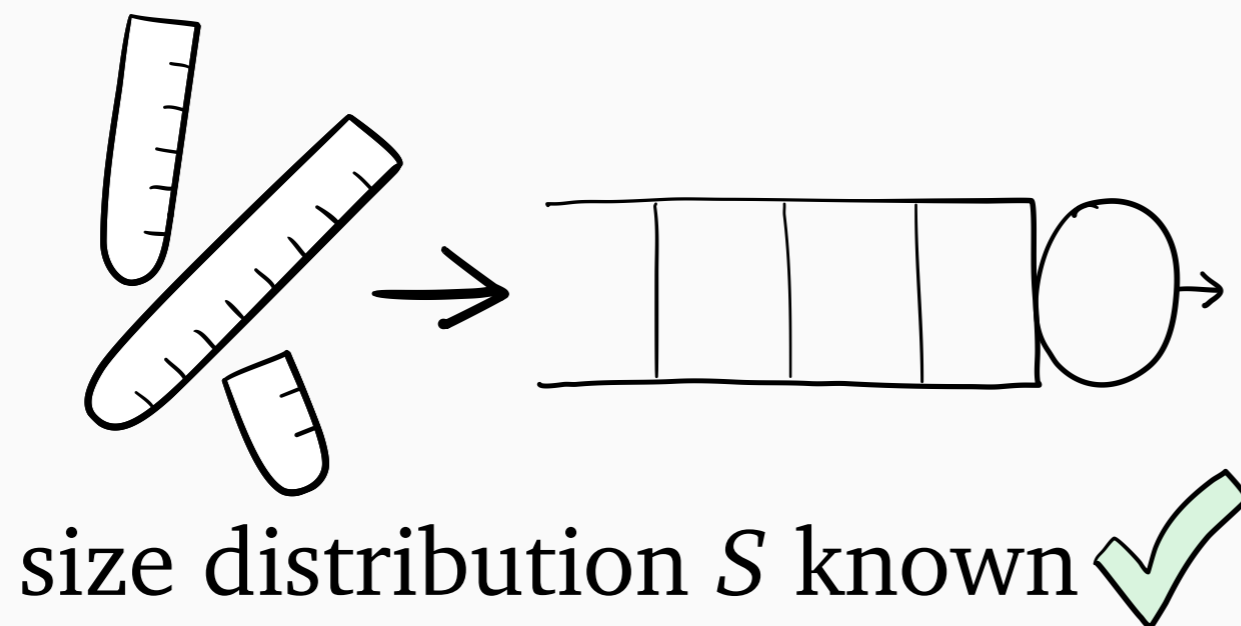
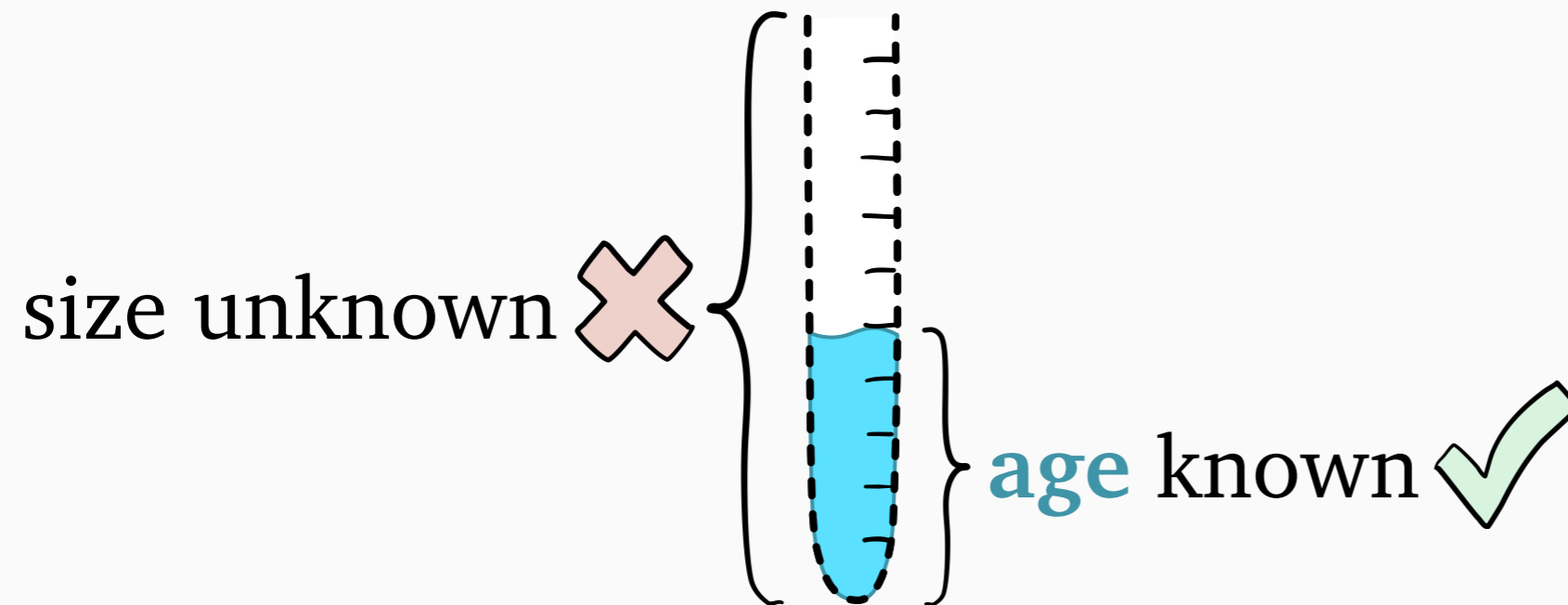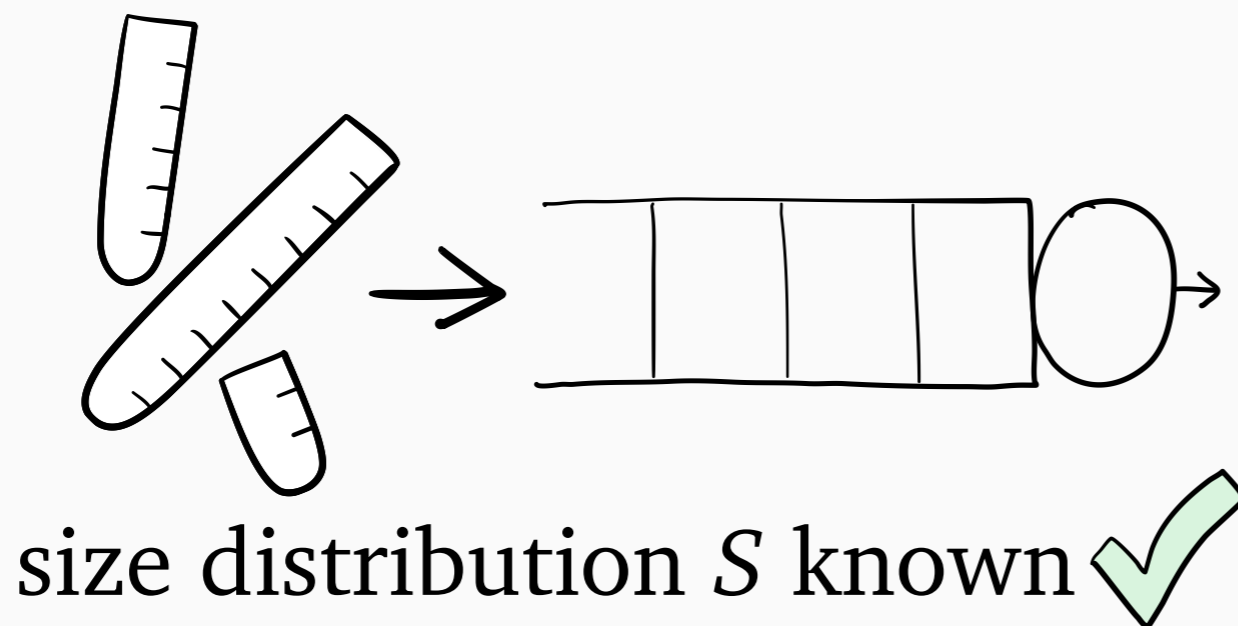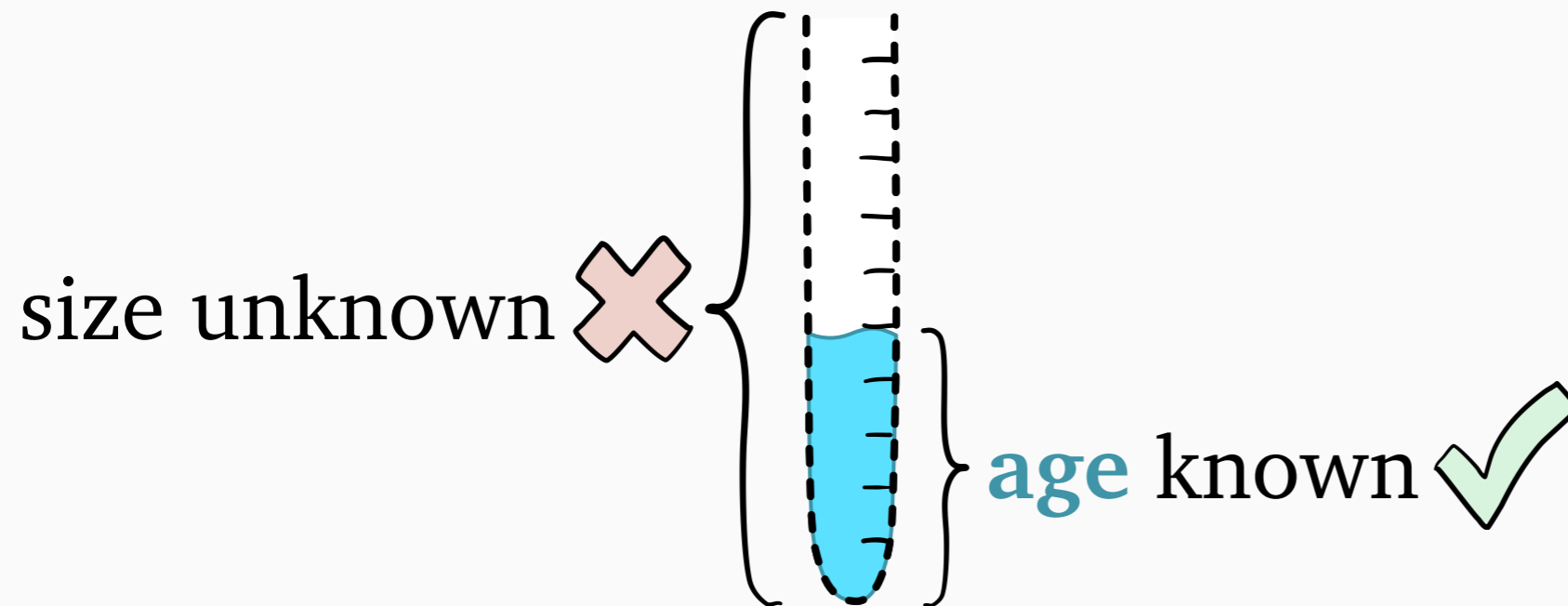# Scheduling with **unknown** sizes

# Scheduling with **unknown** sizes

size unknown ✕

# Scheduling with **unknown** sizes

size unknown ✗

**age** known ✓

# Scheduling with **unknown** sizes



size unknown ✗

**age** known ✓

size distribution $S$ known ✓

# Scheduling with **unknown** sizes



size unknown ✖

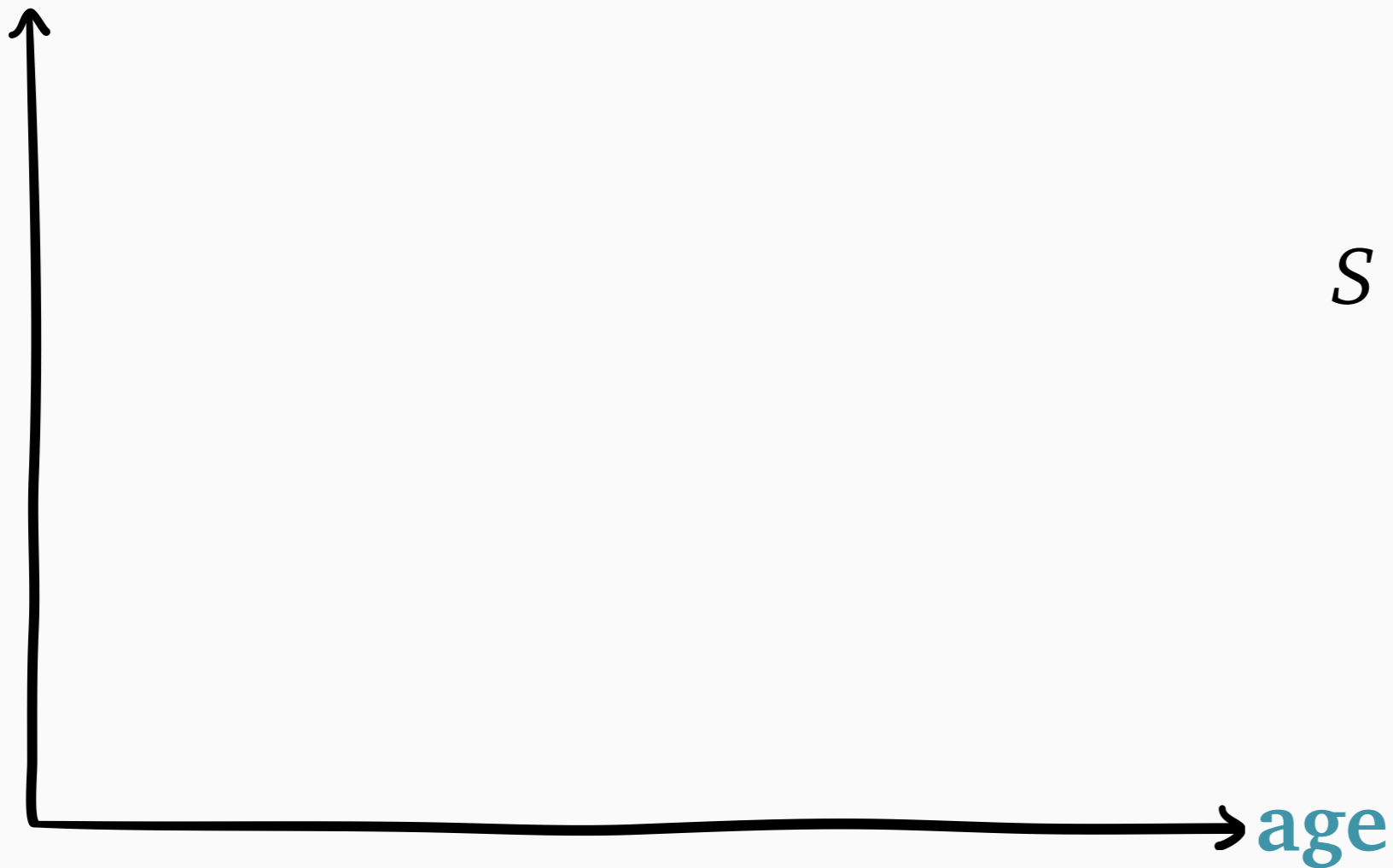**age** known ✔

size distribution $S$ known ✔

**Example:**

$$S = \begin{cases} 1 & \text{w.p.} \ \frac{1}{3} \\ 6 & \text{w.p.} \ \frac{1}{3} \\ 14 & \text{w.p.} \ \frac{1}{3} \end{cases}$$
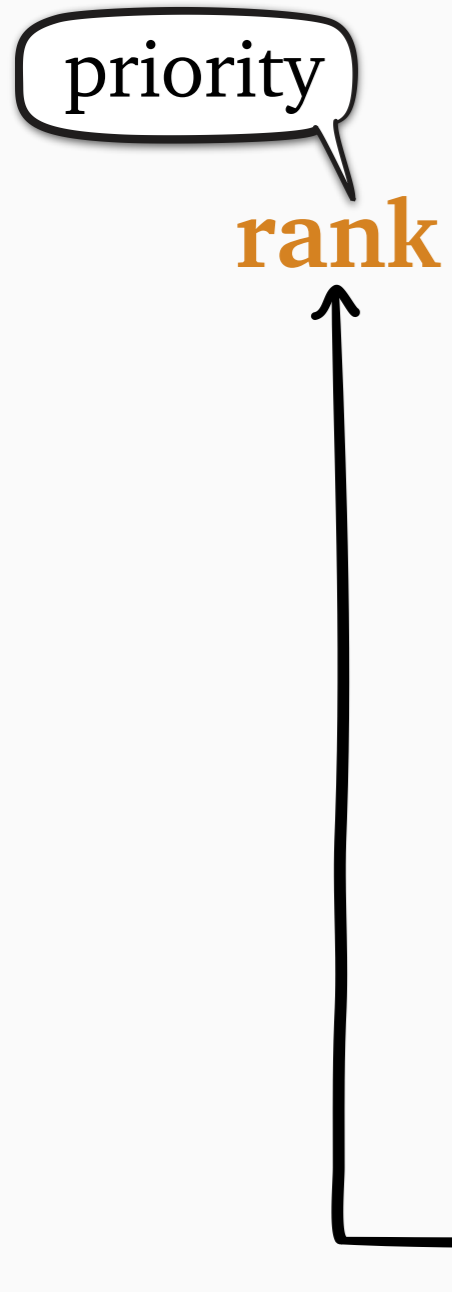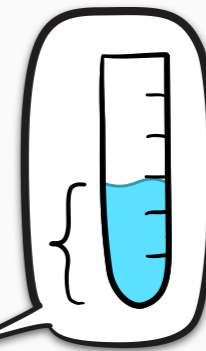
# Scheduling with **unknown** sizes

**rank**

**age**

**Example:**

$$S = \begin{cases} 1 & \text{w.p.} & \frac{1}{3} \\ 6 & \text{w.p.} & \frac{1}{3} \\ 14 & \text{w.p.} & \frac{1}{3} \end{cases}$$
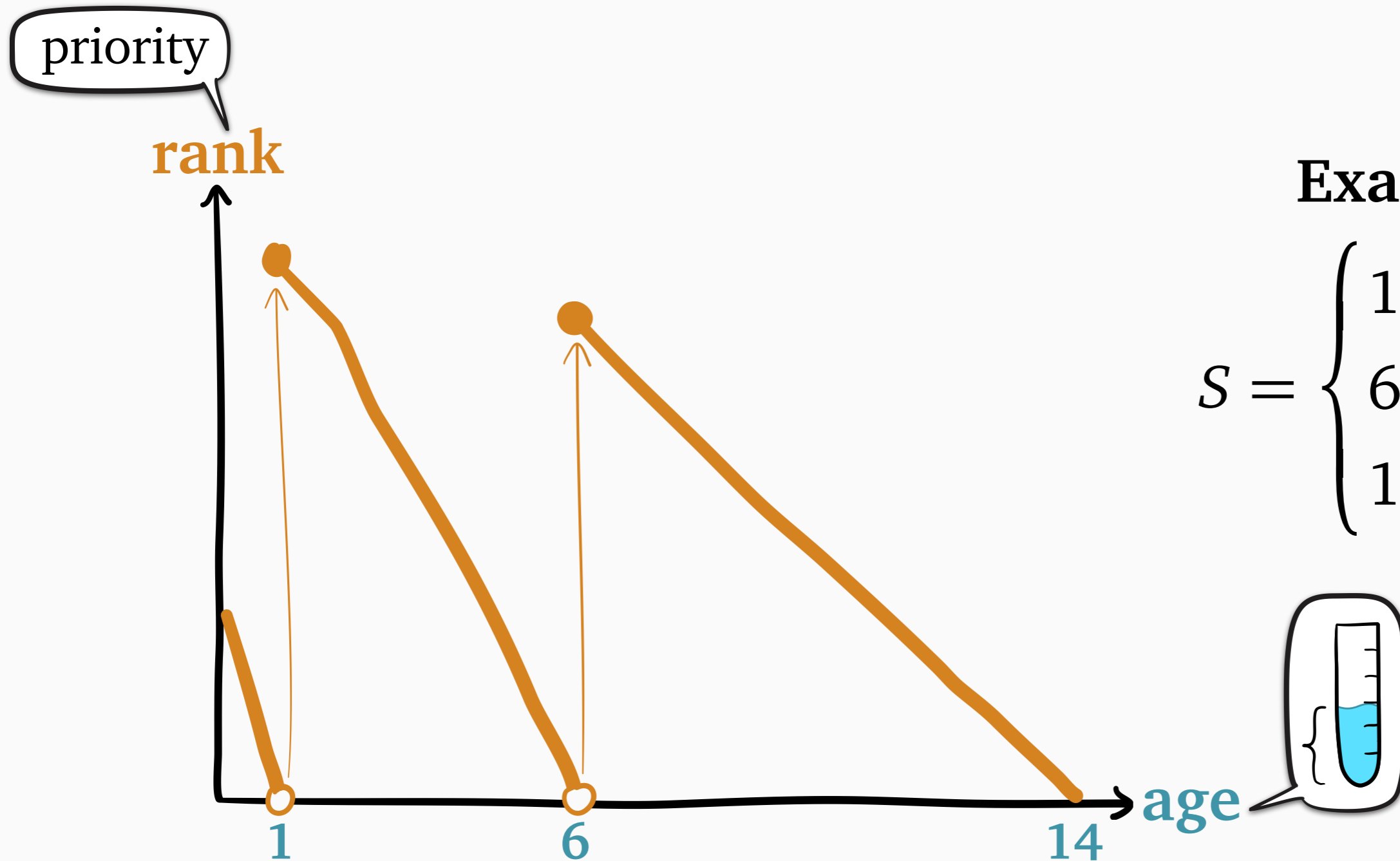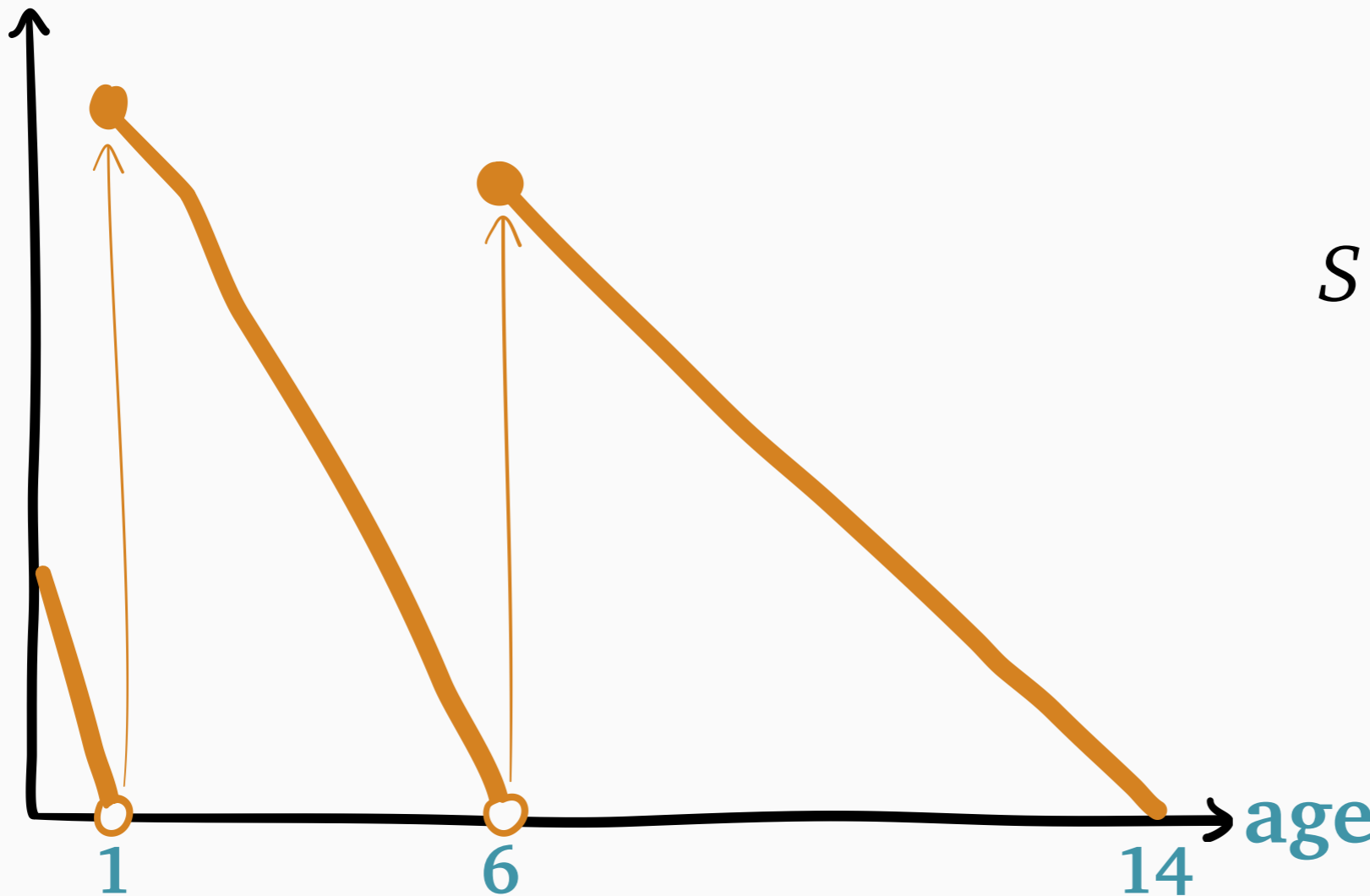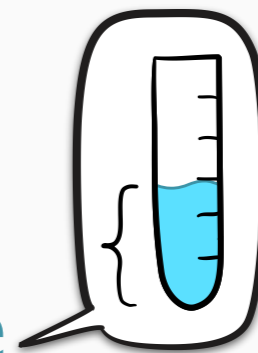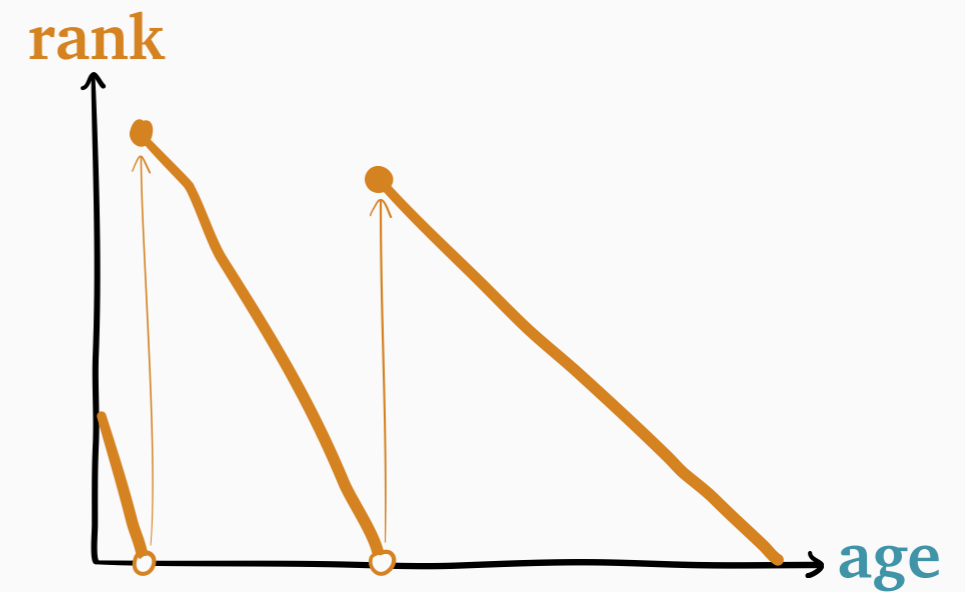
# Scheduling with **unknown** sizes

priority

**rank**

**Example:**

$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$

**age**

# Scheduling with **unknown** sizes



**Example:**

$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$
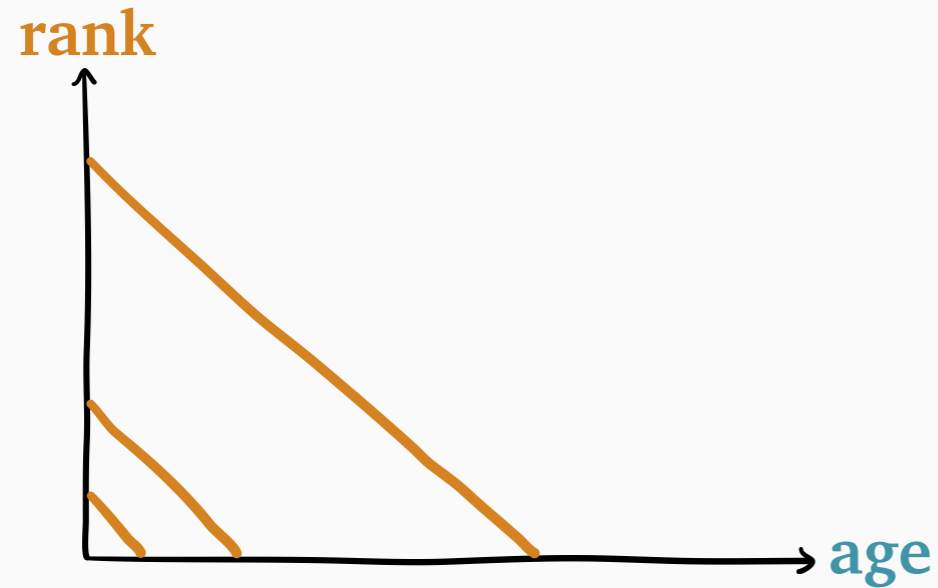
# Scheduling with **unknown** sizes

# Analyzing *nonmonotonic* **rank**



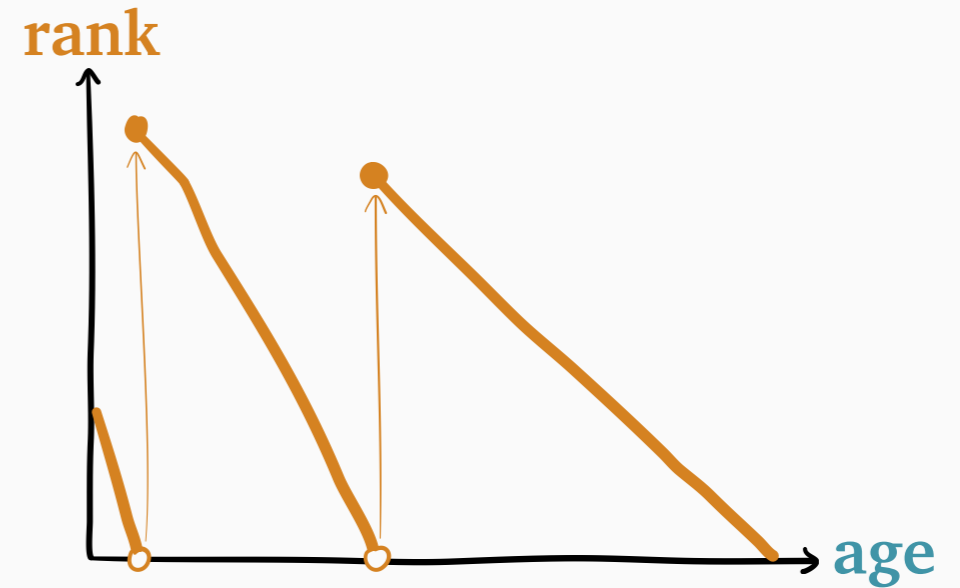**Gittins**: *nonmonotonic*

rank

age

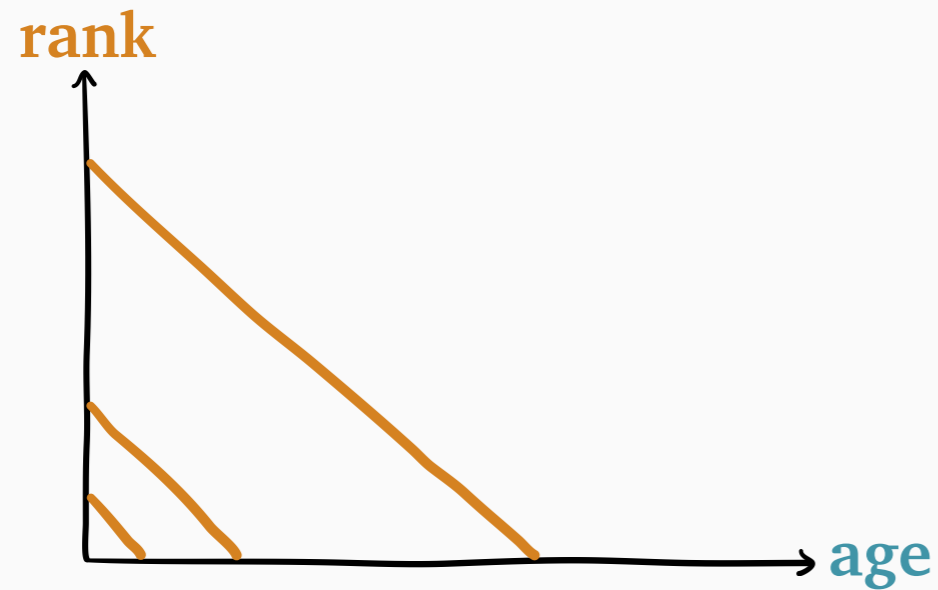# Analyzing *nonmonotonic* **rank**
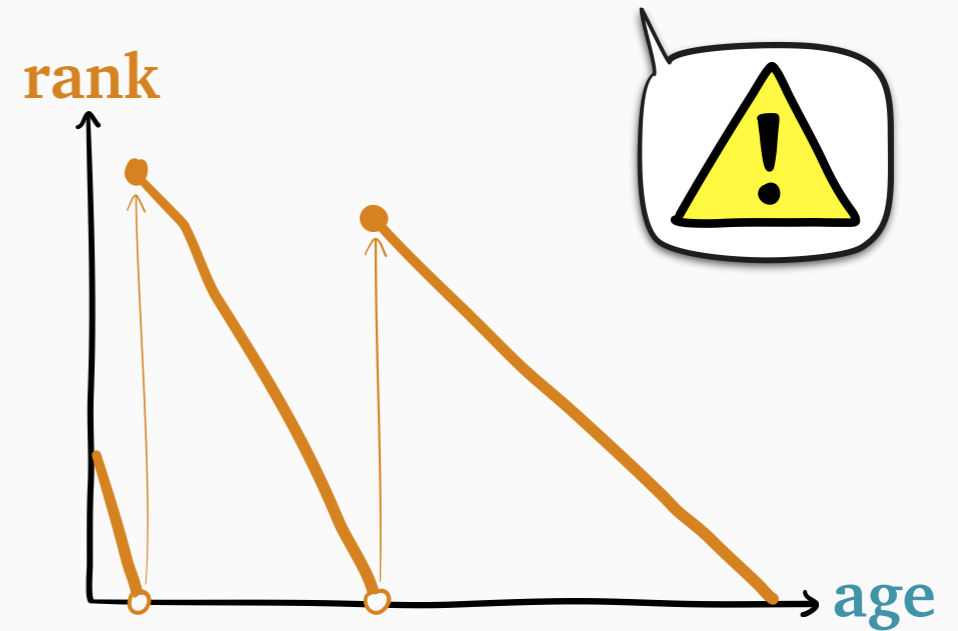
SRPT: monotonic

**Gittins**: *nonmonotonic*

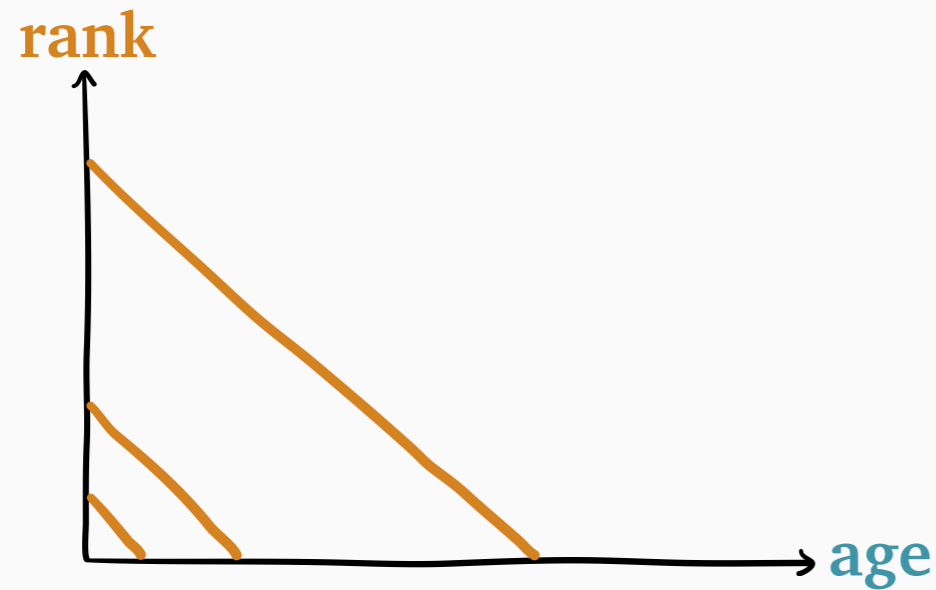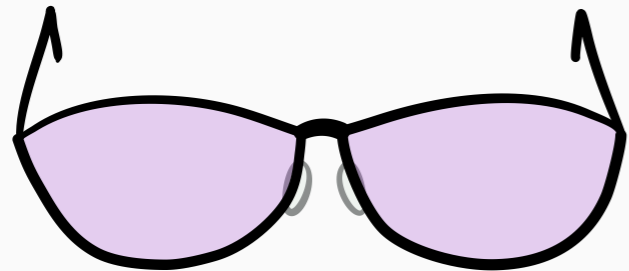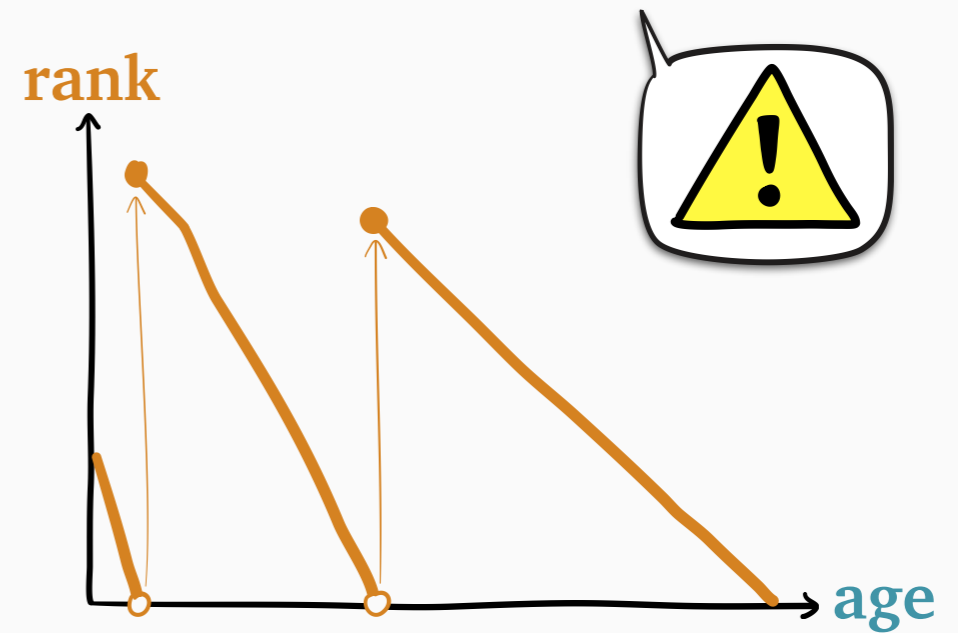# Analyzing *nonmonotonic* **rank**
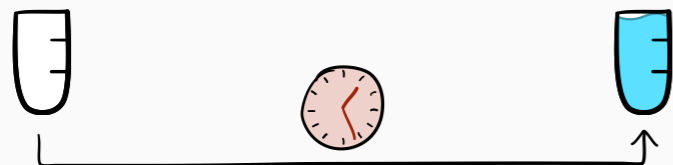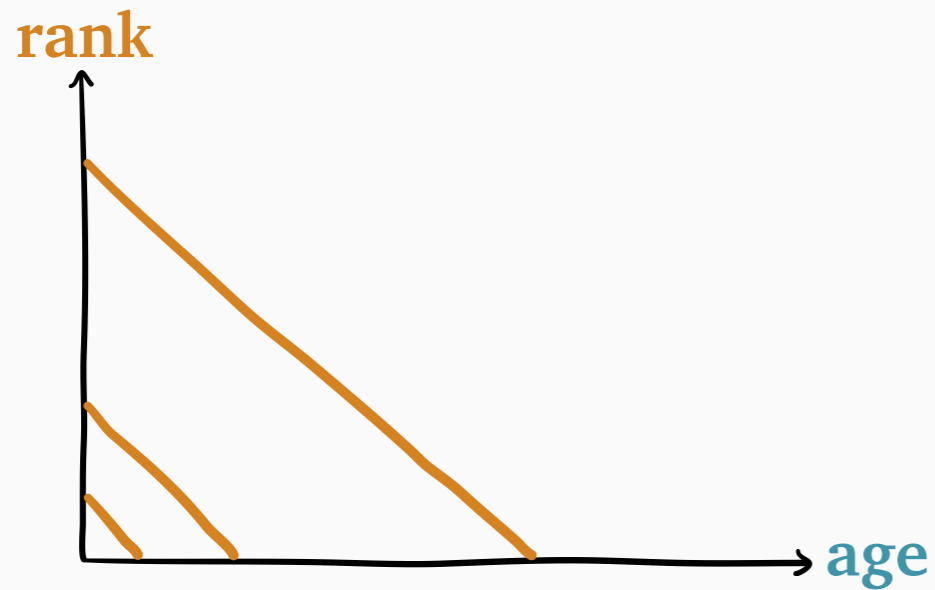
SRPT: monotonic

Gittins: *nonmonotonic*

# Analyzing *nonmonotonic* **rank**

SRPT: monotonic

**Gittins**: *nonmonotonic*



**Step 1:** compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$
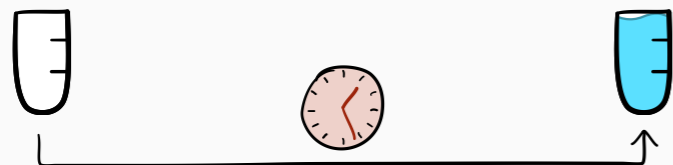
# Analyzing *nonmonotonic* **rank**

SRPT: monotonic

**Gittins**: *nonmonotonic*



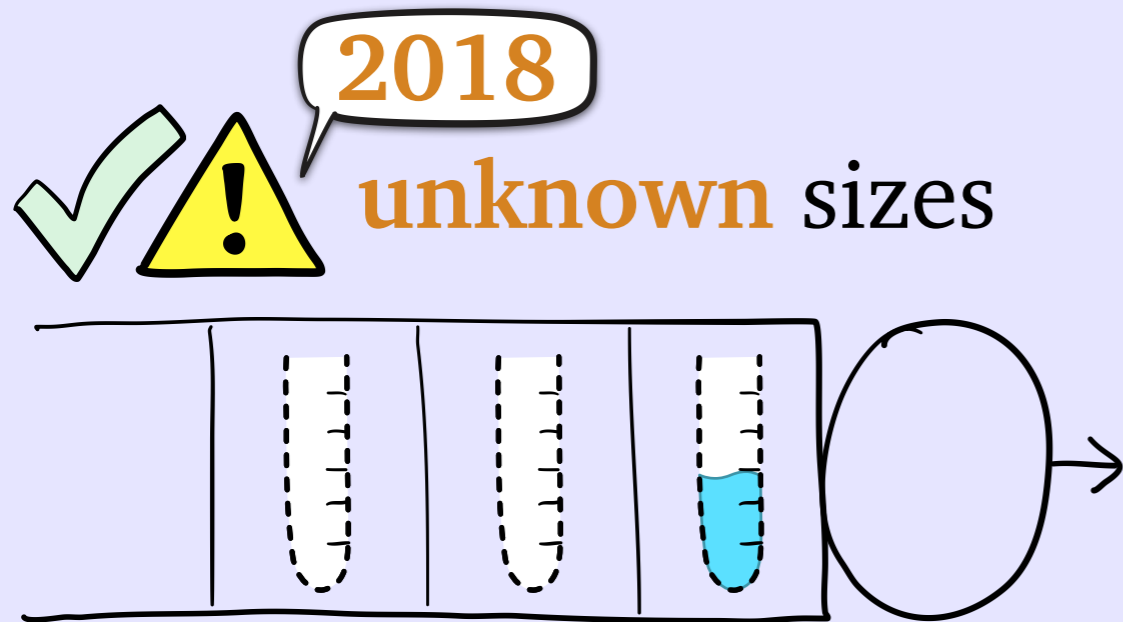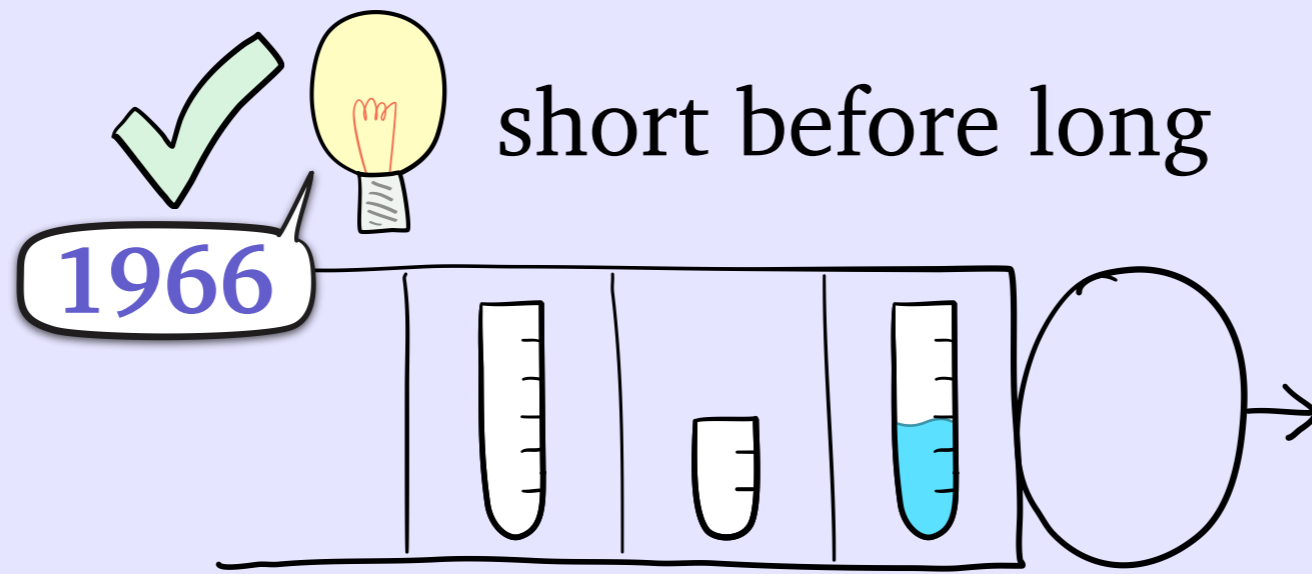**Step 1:** compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

*fancy*
**tagged job**

Tagged job methods

short before long — 1966 ✓

unknown sizes — 2018 ✓⚠

multiple servers — 2018 ✓⚠

This work: both at once! — 2021 ⚠⚠

13

**unknown** sizes + **multiple** servers

**unknown** sizes + **multiple** servers



**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$

*k* servers: *intractable*

**unknown** sizes + **multiple** servers

**Step 1:** compute $\mathbf{E}[W(r)]$

$k$ servers:
*intractable*

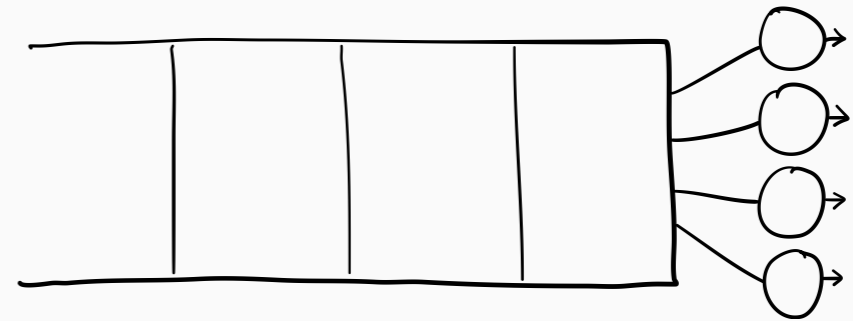**unknown** sizes + **multiple** servers

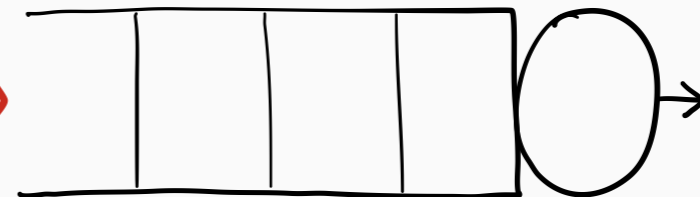**Step 1:** compute $\mathbf{E}[W(r)]$

SRPT's $W(r)$ gap $\leq kr$

k servers: *intractable*

**unknown** sizes + **multiple** servers

**Step 1:** compute $\mathbf{E}[W(r)]$
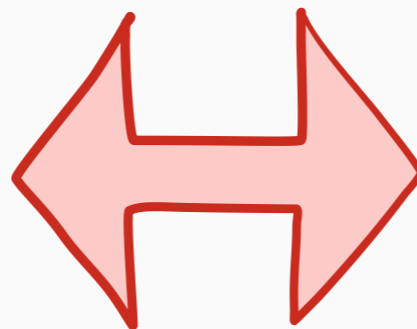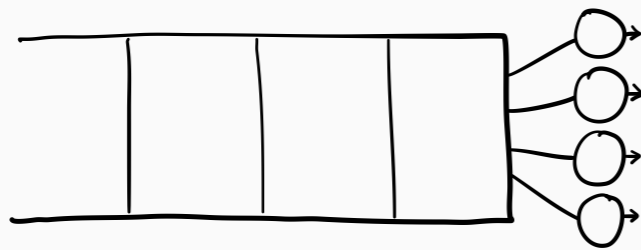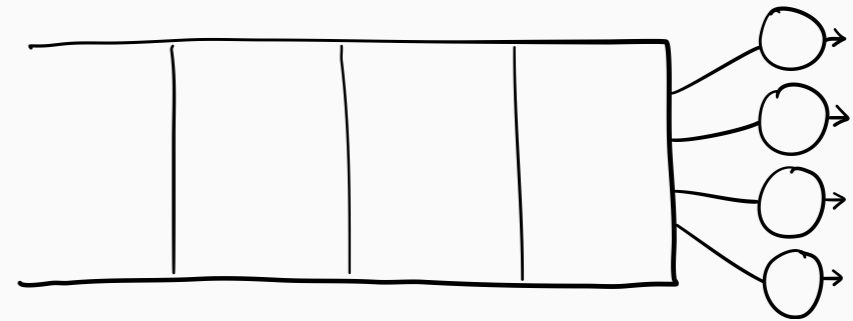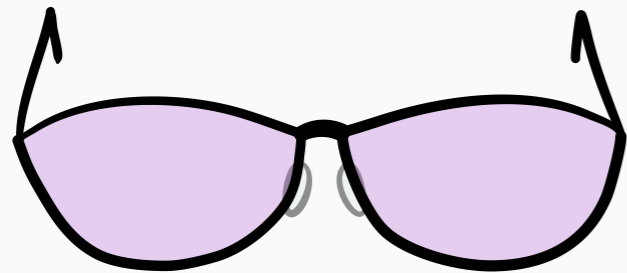
SRPT's $W(r)$ gap $\leq kr$
**Gittins**'s $W(r)$ gap $\leq k\infty$

**unknown** sizes + **multiple** servers

*k* servers: *intractable*

**Step 1:** compute $\mathbf{E}[W(r)]$

SRPT's $W(r)$ gap $\leq kr$

Gittins's $W(r)$ gap $\leq k\infty$

unknown + worst-case

**unknown** sizes + **multiple** servers

Step 1: compute $\mathbf{E}[W(r)]$
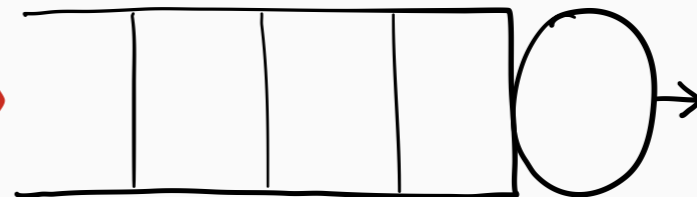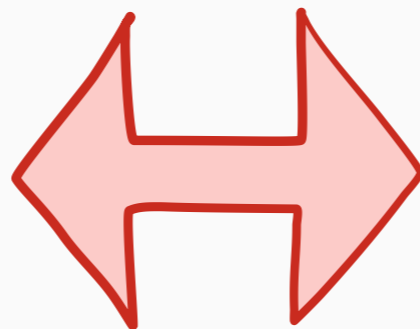
**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$
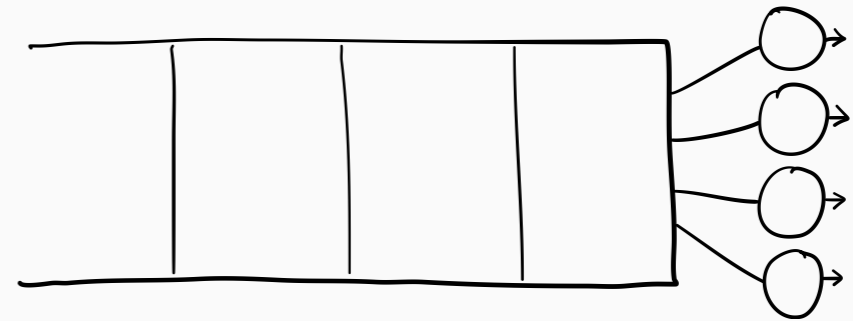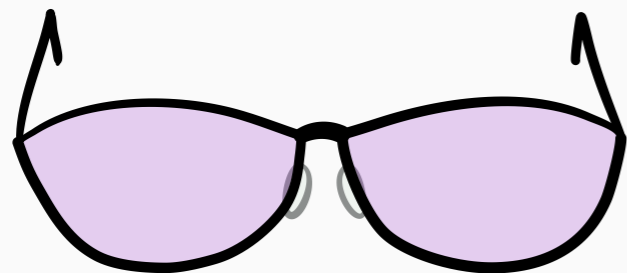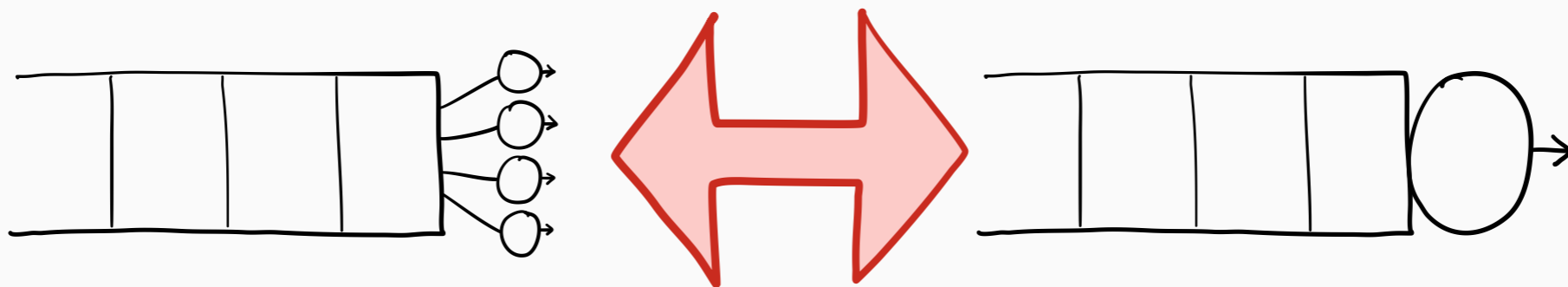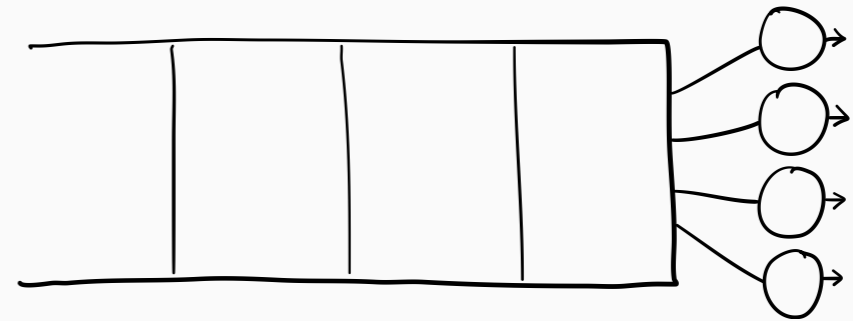
**$k$** servers: *intractable*

**unknown** sizes + **multiple** servers

Step 1: compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

**tagged job**
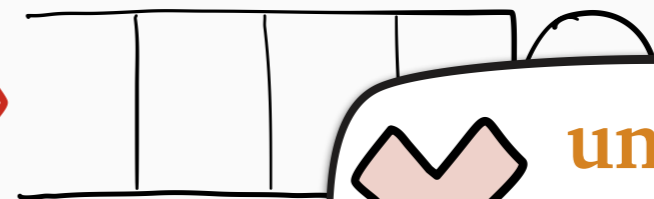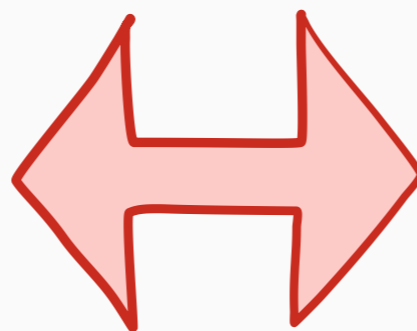**+**
**worst-case**

**Gittins rank**: *nonmonotonic*

*k* servers: *intractable*

**unknown** sizes + **multiple** servers

Step 1: compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

**tagged job**
**+**
**worst-case**

**Gittins rank**: *nonmonotonic*

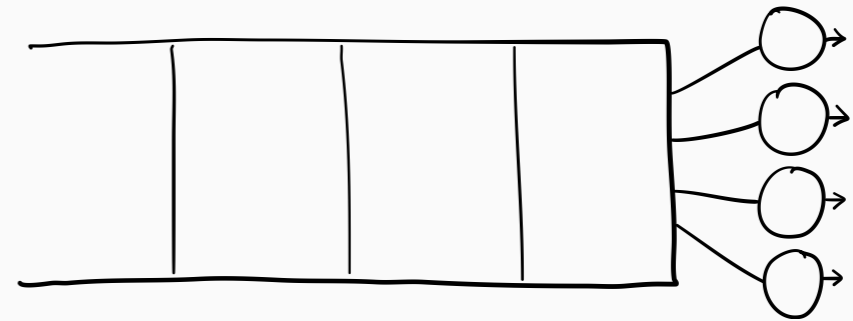*k* servers: *intractable*

**unknown** sizes + **multiple** servers

Step 1: compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

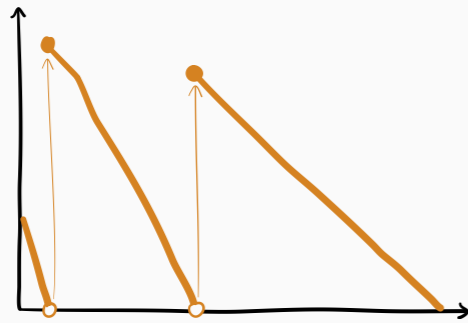*fancy* **tagged job**
+
**worst-case**

**Gittins rank**: *nonmonotonic*

*k* servers: *intractable*

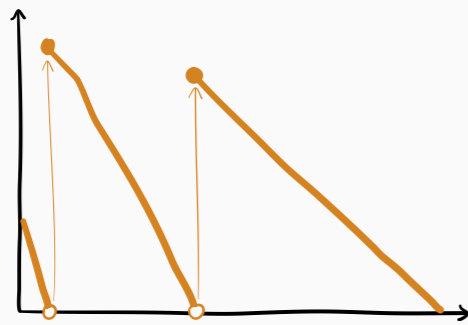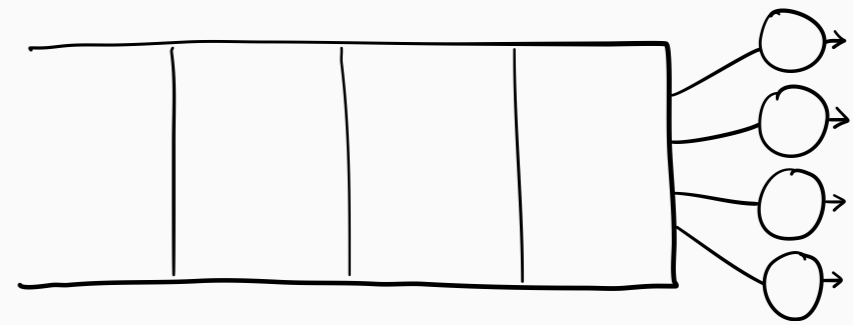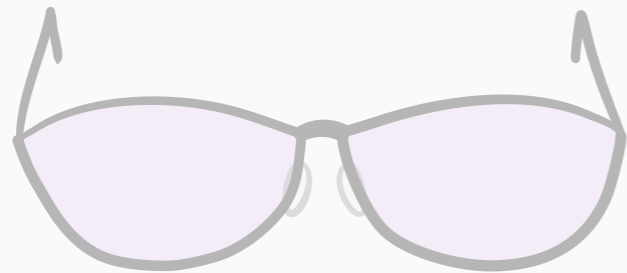**unknown** sizes + **multiple** servers

Step 1: compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

*fancy* **tagged job**
+
**worst-case**

**Step 1:** compute $\mathbb{E}[W(r)]$

**Step 2:** $\mathbb{E}[W(r)]$ to $\mathbb{E}[T]$

**Step 1:** compute $\mathbf{E}[W(r)]$
*without **worst-case** steps*

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$
*without **worst-case** steps*

**Step 2:** $\mathbf{E}[W(\textcolor{purple}{r})]$ to $\mathbf{E}[T]$
*without **tagged job** method*

**Step 1:** compute $\mathbf{E}[W(r)]$
*without **worst-case** steps*

**Idea:**

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$
*without **tagged job** method*

**Step 1:** compute $\mathbf{E}[W(\textcolor{purple}{r})]$
*without **worst-case** steps*

**Idea:**

**Step 2:** $\mathbf{E}[W(\textcolor{purple}{r})]$ to $\mathbf{E}[T]$
*without **tagged job** method*

**Idea:** $s$

**Step 1:** compute $\mathbf{E}[W(\textbf{\textit{r}})]$
*without **worst-case** steps*

**Idea:**

**Step 2:** $\mathbf{E}[W(\textbf{\textit{r}})]$ to $\mathbf{E}[T]$
*without **tagged job** method*

**Idea:** $s$

# Our contribution:
new **exact formulas** for both steps

**Step 1:** compute $\mathbf{E}[W(r)]$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

**Step 1:** compute $\mathbf{E}[W(r)]$



## Work Decomposition Law

$$\mathbf{E}[W_k(r)] = \mathbf{E}[W_1(r)] + \mathbf{E}[\text{``} < k \text{ jobs' } r\text{-work''}]$$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

**Step 1:** compute $\mathbf{E}[W(r)]$



**Work Decomposition Law**

$$\mathbf{E}[W_k(r)] = \mathbf{E}[W_1(r)] + \mathbf{E}[``< k \text{ jobs' } r\text{-work''}]$$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$



**WINE**

$$\lambda \mathbf{E}[T] = \mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(r)]}{r^2}\, \mathrm{d}r$$

**Step 1:** compute $\mathbf{E}[W(r)]$



## Work Decomposition Law

$$\mathbf{E}[W_k(r)] = \mathbf{E}[W_1(r)] + \mathbf{E}[\text{``}< k \text{ jobs' } r\text{-work''}]$$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

## WINE



$$\lambda\mathbf{E}[T] = \mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(r)]}{r^2}\,\mathrm{d}r$$

specific to **Gittins**'s **rank**

16

**Step 1:** compute $\mathbf{E}[W(r)]$

## Work Decomposition Law

$$\mathbf{E}[W_k(r)] = \mathbf{E}[W_1(r)] + \mathbf{E}[\text{``} < k \text{ jobs' } r\text{-work''}]$$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

## WINE

specific to **Gittins**'s **rank**

$$\lambda\mathbf{E}[T] = \mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(r)]}{r^2}\,\mathrm{d}r$$

**Impact 1:**
first bound on $\mathbf{E}[T_{\text{Gittins-}k}]$

**Step 1:** compute $\mathbf{E}[W(r)]$



**Work Decomposition Law**

$$\mathbf{E}[W_k(r)] = \mathbf{E}[W_1(r)] + \mathbf{E}[\text{“}< k \text{ jobs' } r\text{-work”}]$$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$



**WINE**

$$\lambda\mathbf{E}[T] = \mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(r)]}{r^2}\,\mathrm{d}r$$

specific to **Gittins**'s **rank**

**Impact 1:**
first bound on $\mathbf{E}[T_{\text{Gittins-}k}]$

**Impact 2:**
both generalize beyond M/G/$k$

**Step 1:** compute $\mathbf{E}[W(r)]$



**Work Decomposition Law**

$$\mathbf{E}[W_k(r)] = \mathbf{E}[W_1(r)] + \mathbf{E}[\text{``}< k \text{ jobs' } r\text{-work''}]$$

**Step 2:** $\mathbf{E}[W(r)]$ to $\mathbf{E}[T]$

**WINE**

$$\lambda \mathbf{E} \qquad \int_0^\infty \frac{\mathbf{E}[W(r)]}{\phantom{x}} \, dr$$

specific to **Gittins**'s **rank**

**Example:** load-balancing + scheduling

**Impact 1:**
first bound on $\mathbf{E}[T_{\text{Gittins-}k}]$

**Impact 2:**
both generalize beyond M/G/$k$

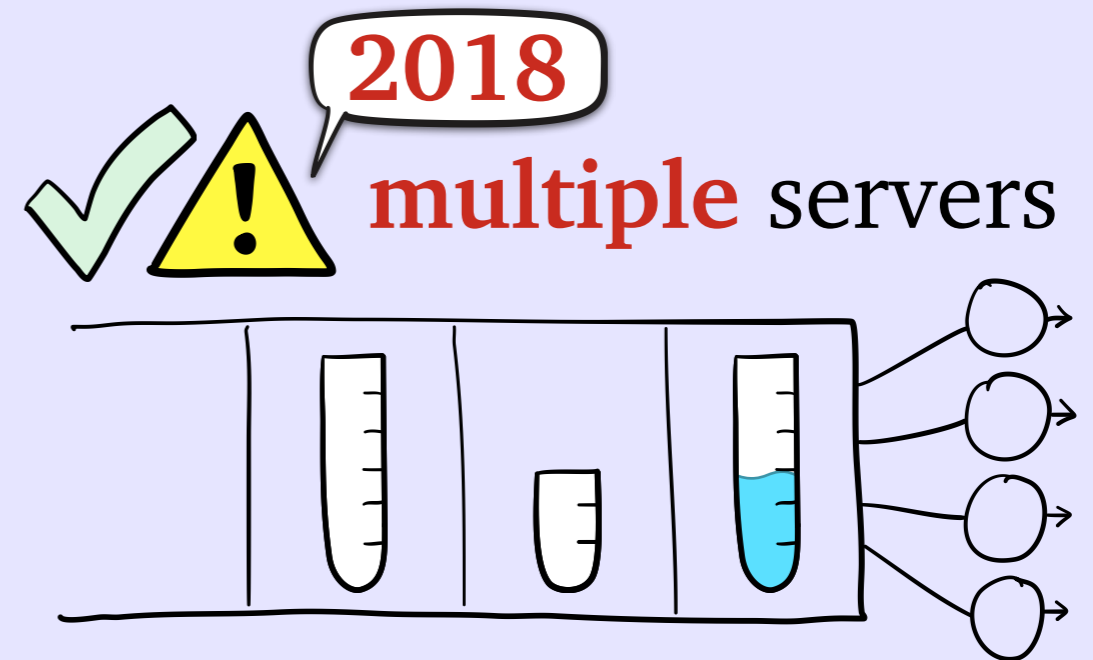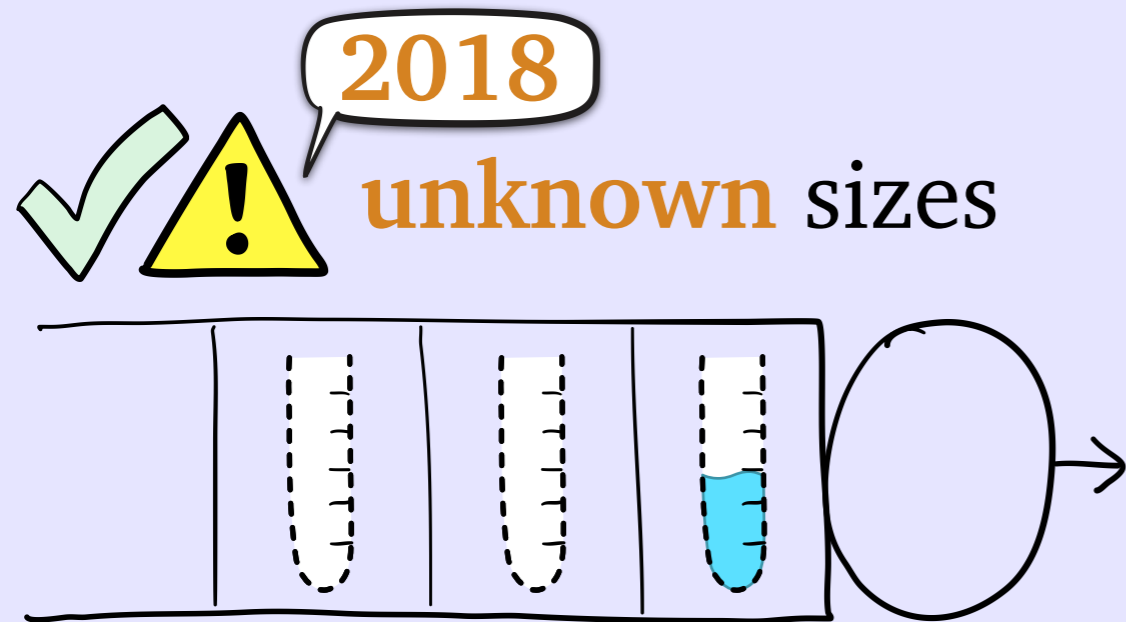16