

What **TCS** Can Do
for **Queueing**



and

What **Queueing**
Can Do for **TCS**

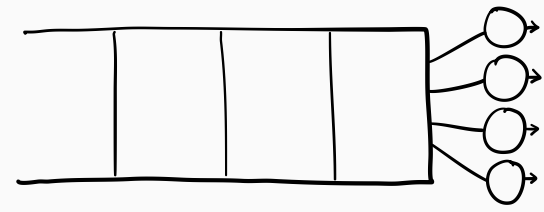
in Scheduling Theory

Ziv Scully
Cornell University



Part I

Handling job size uncertainty



Part II

Analyzing multiserver scheduling

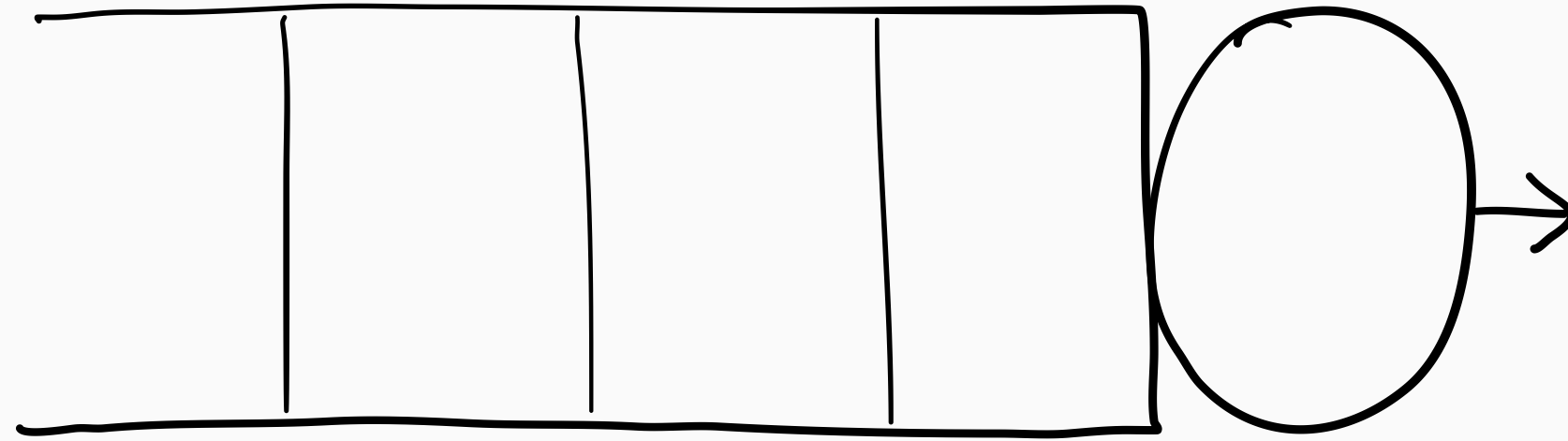


Part III

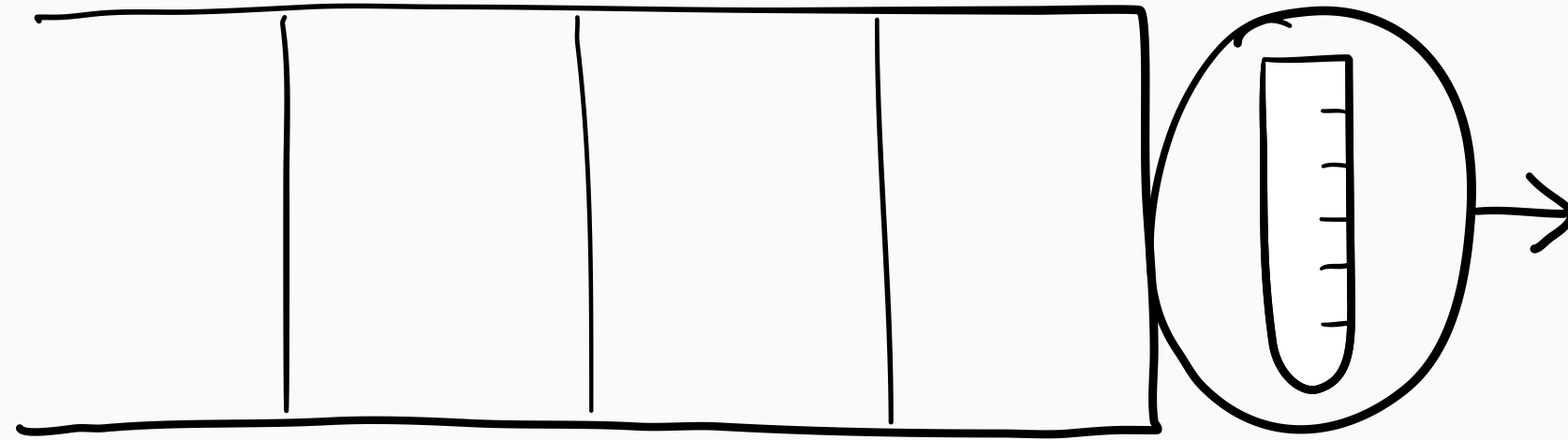
Optimizing tail metrics

How should we schedule jobs to minimize delay?

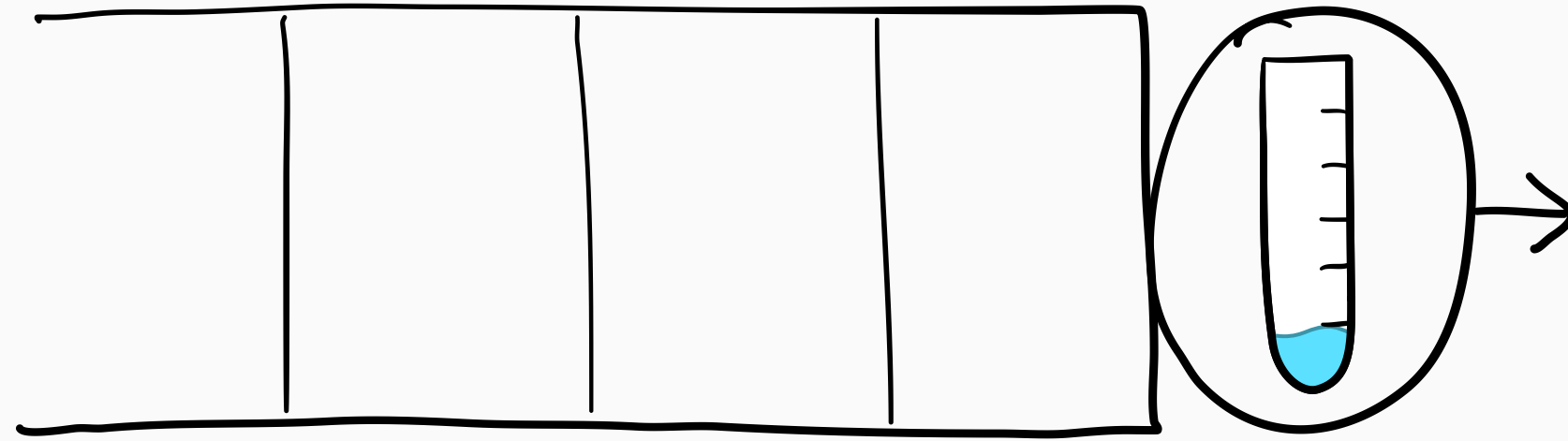
How should we schedule jobs to minimize delay?



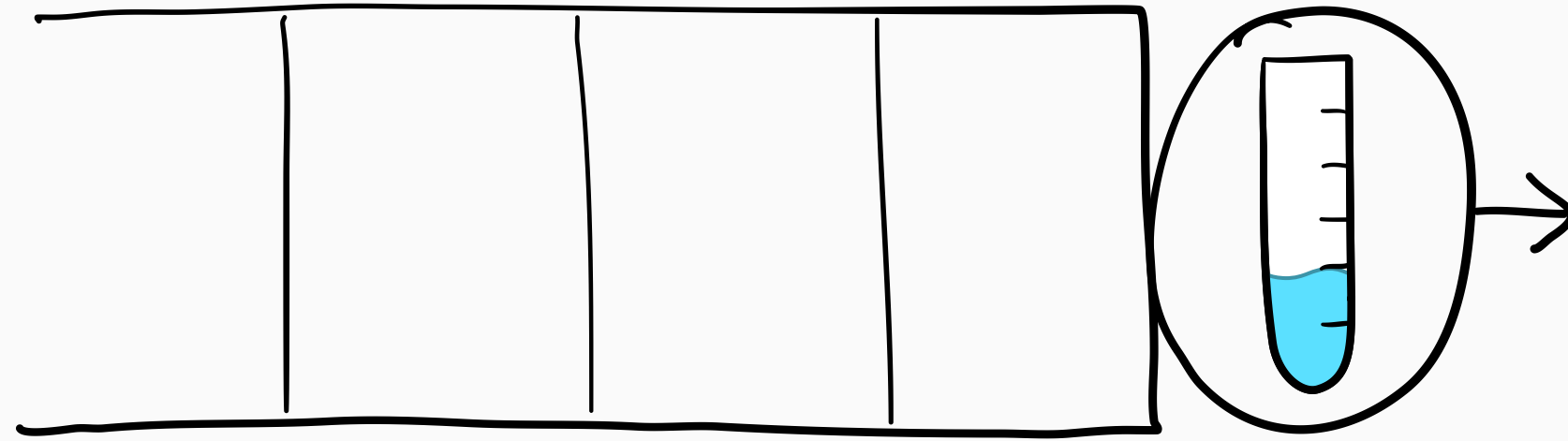
How should we schedule jobs to minimize delay?



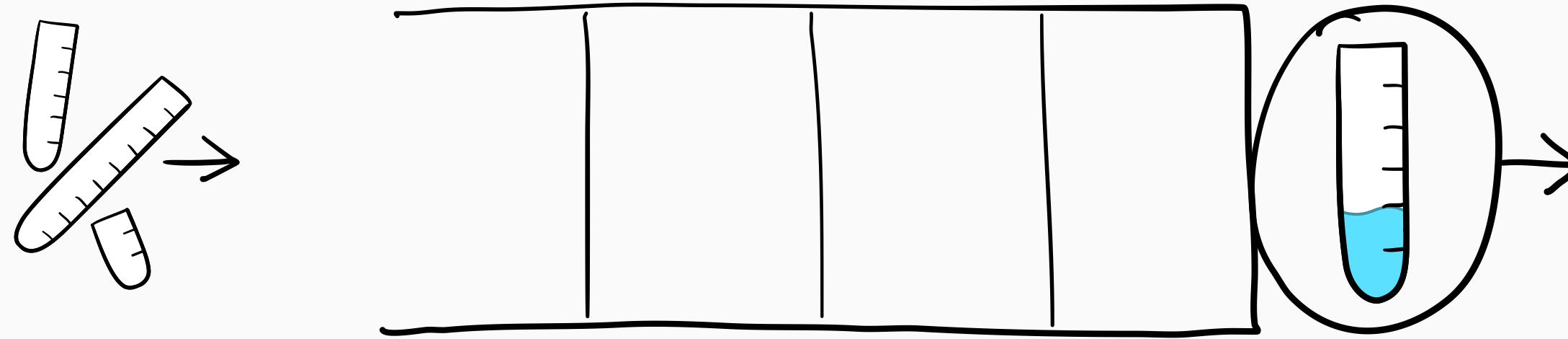
How should we schedule jobs to minimize delay?



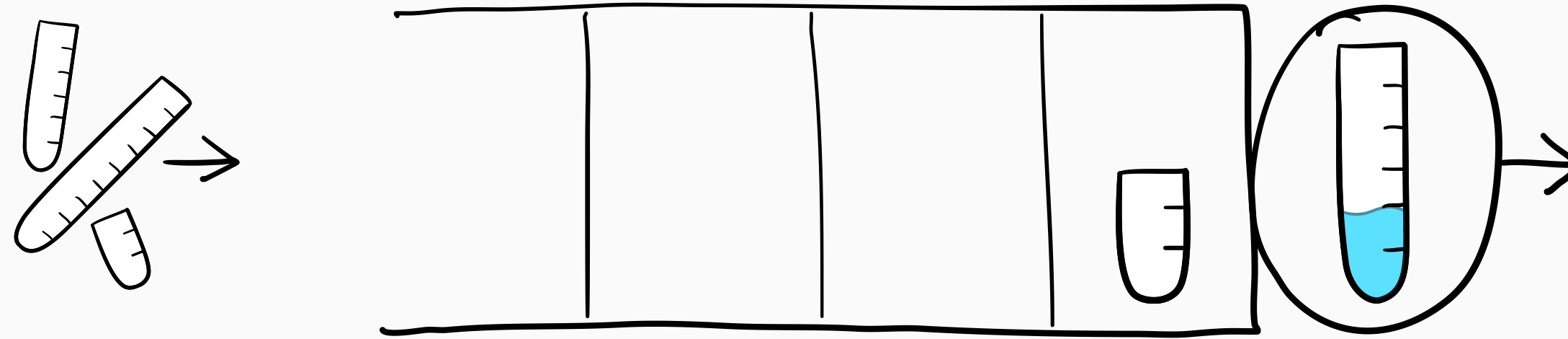
How should we schedule jobs to minimize delay?



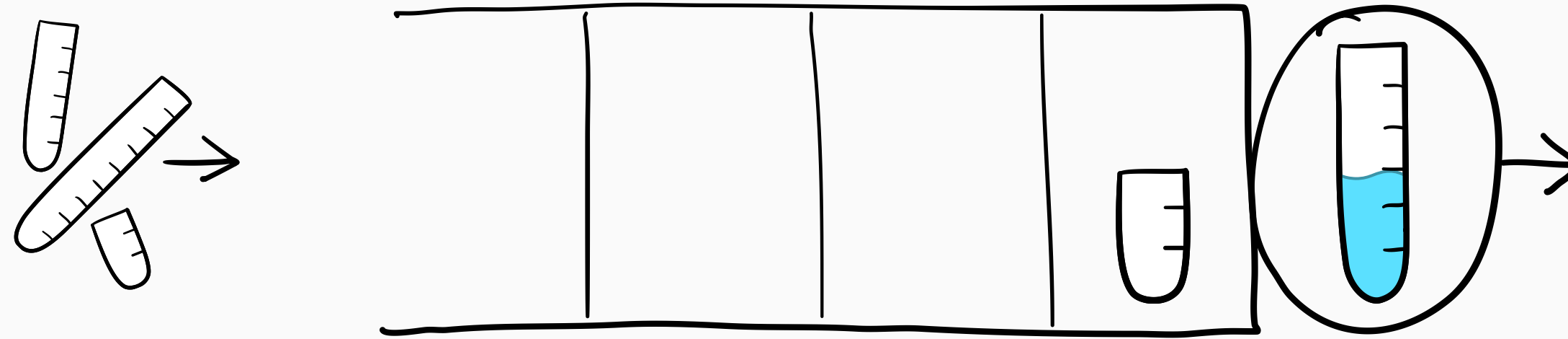
How should we schedule jobs to minimize delay?



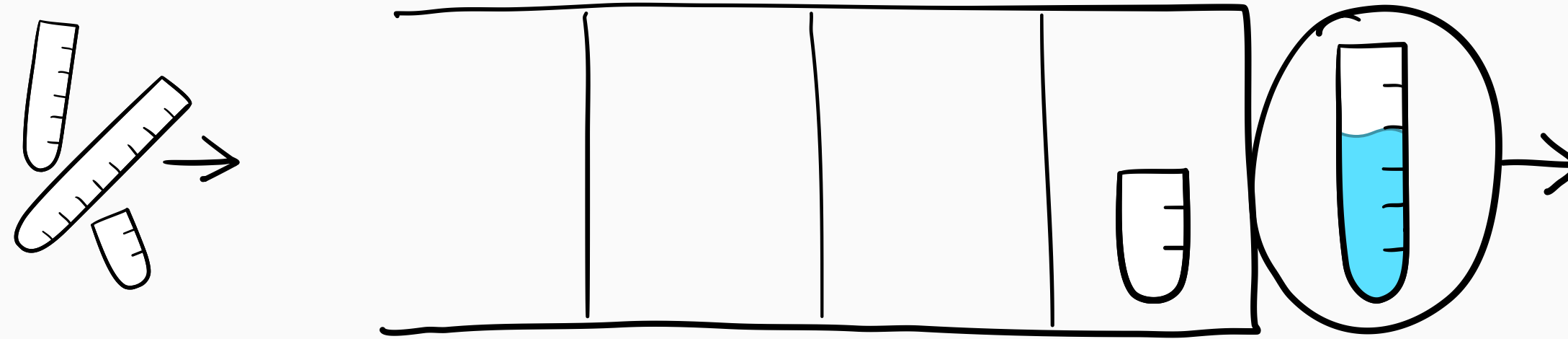
How should we schedule jobs to minimize delay?



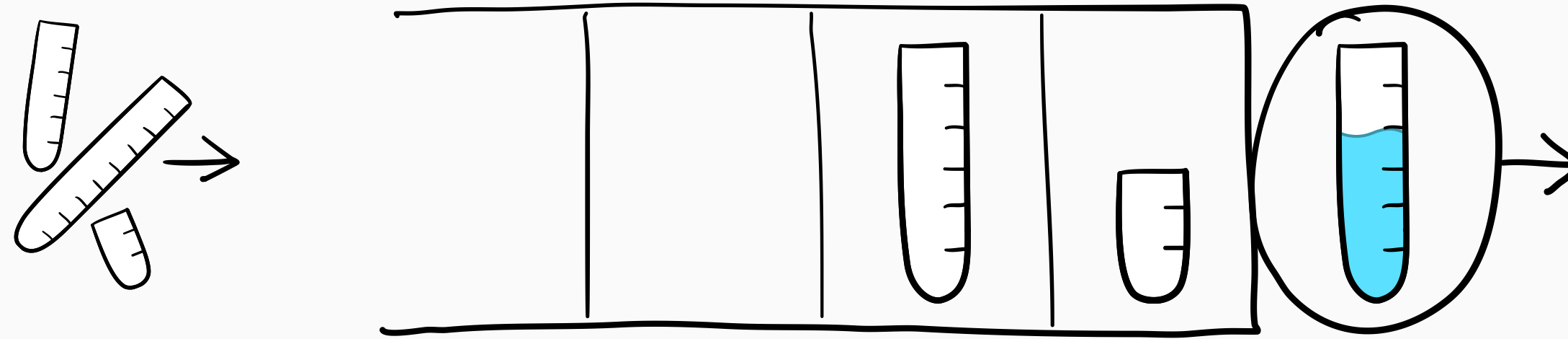
How should we schedule jobs to minimize delay?



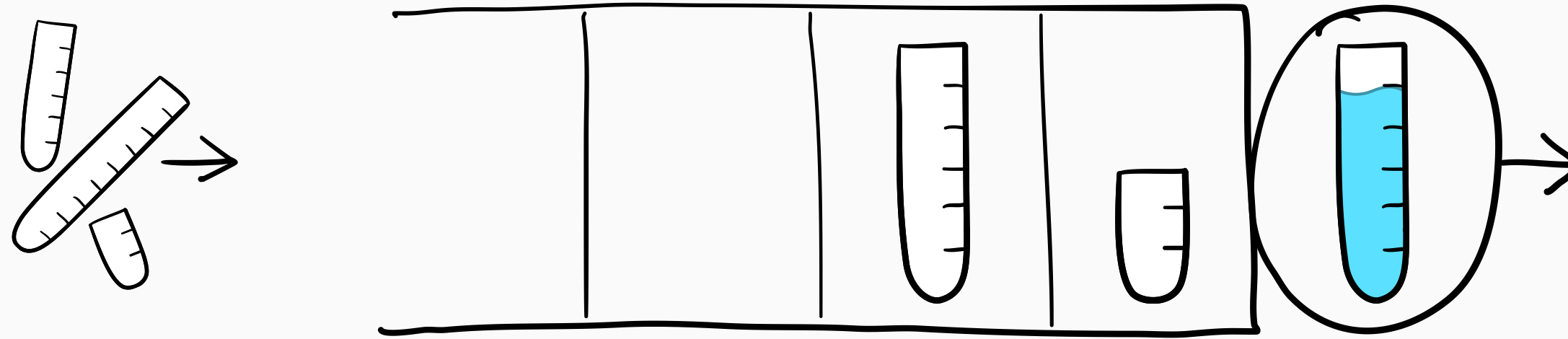
How should we schedule jobs to minimize delay?



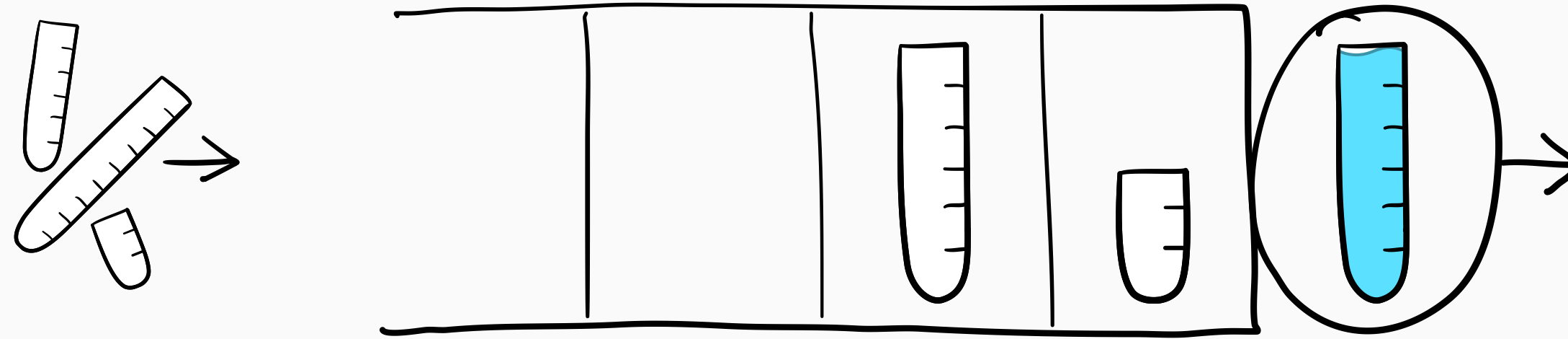
How should we schedule jobs to minimize delay?



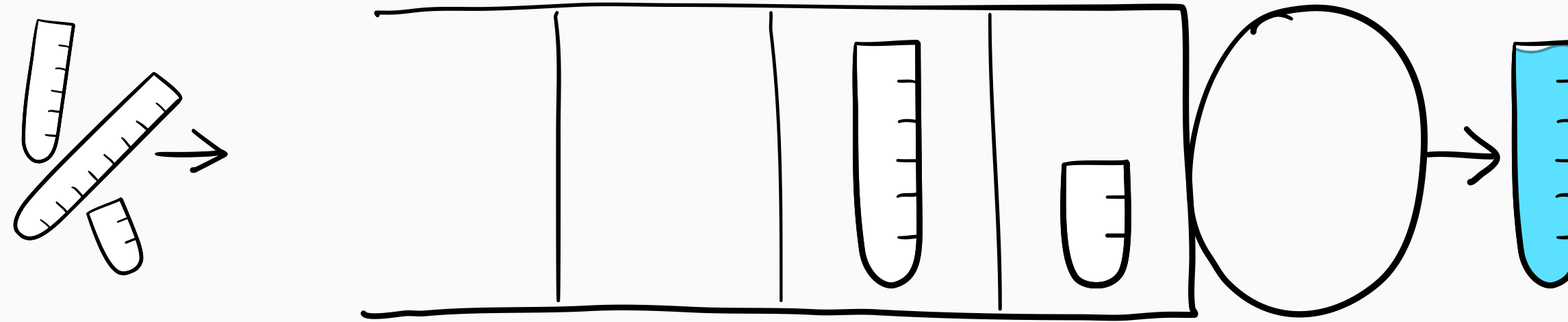
How should we schedule jobs to minimize delay?



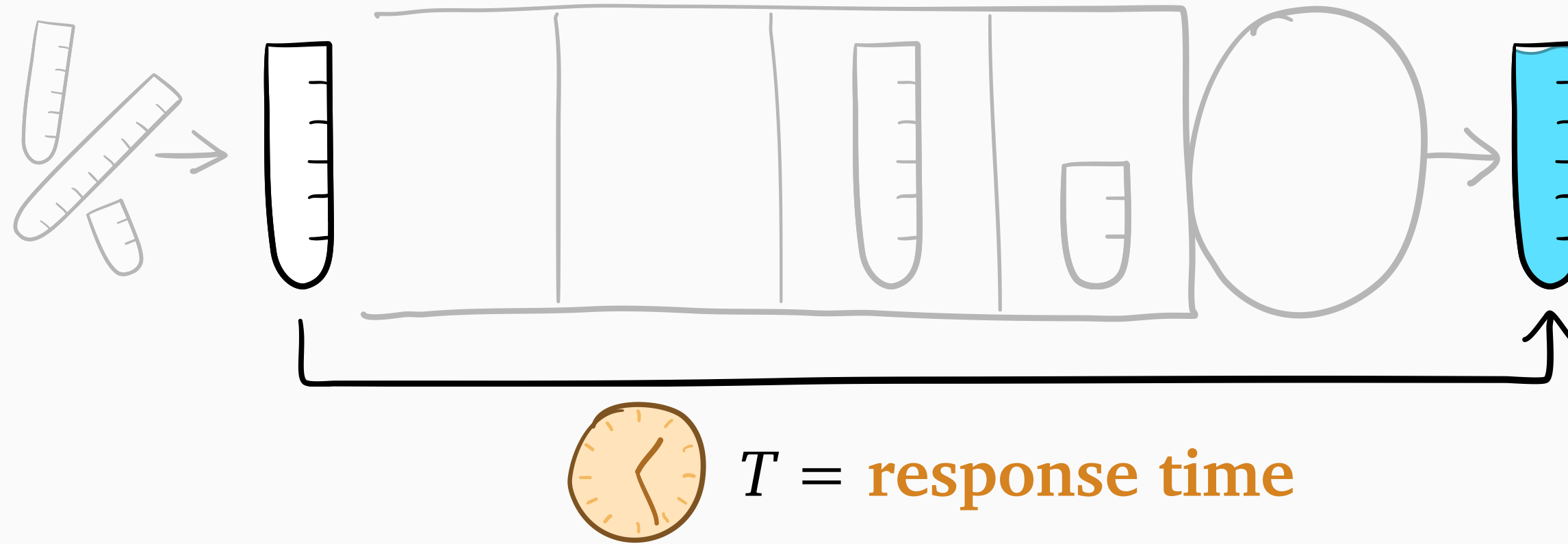
How should we schedule jobs to minimize delay?



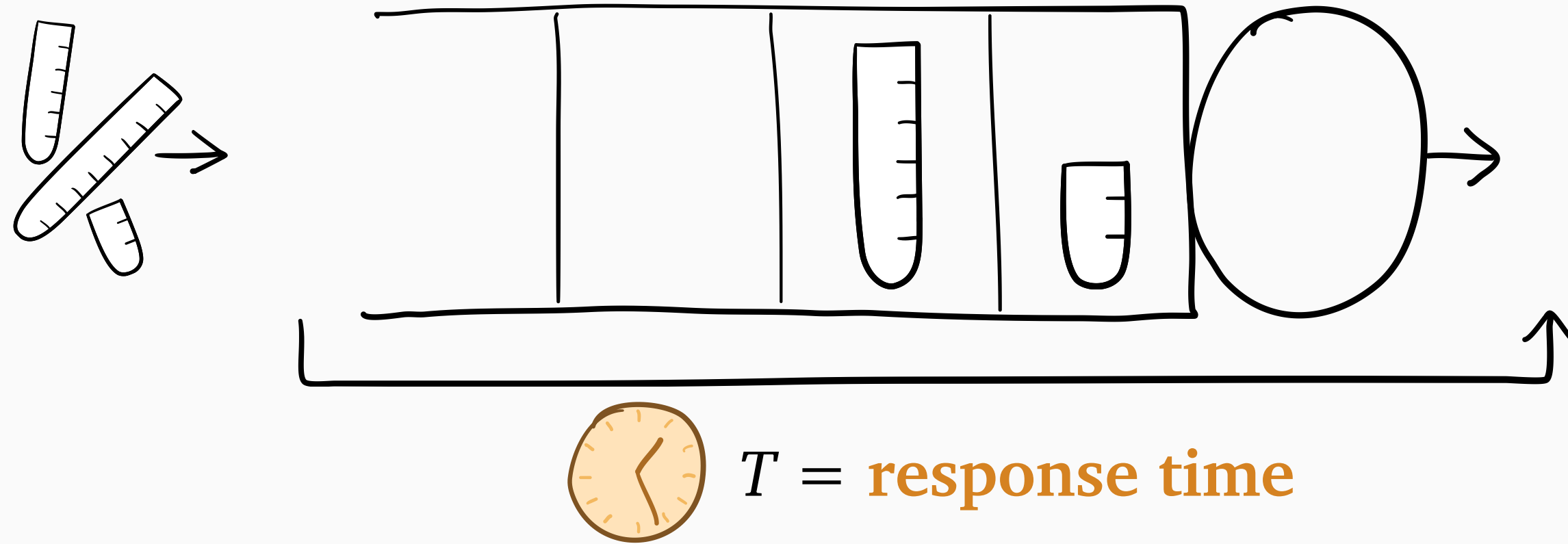
How should we schedule jobs to minimize delay?



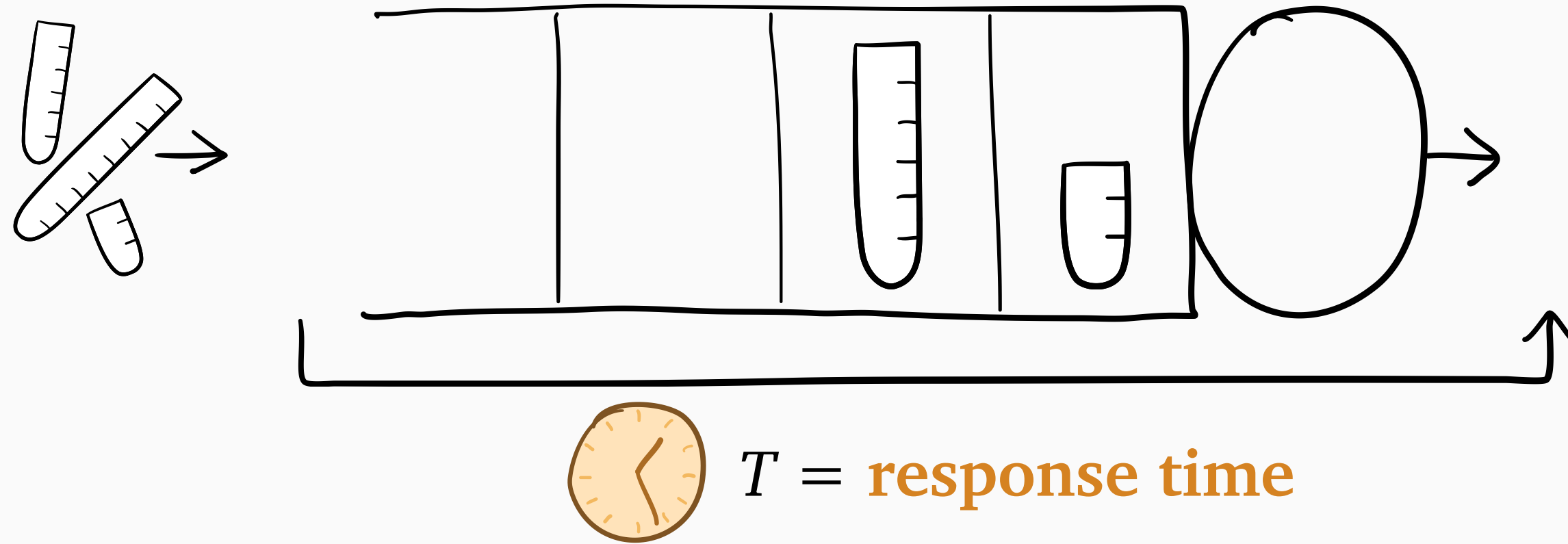
How should we schedule jobs to minimize delay?



How should we schedule jobs to minimize delay?

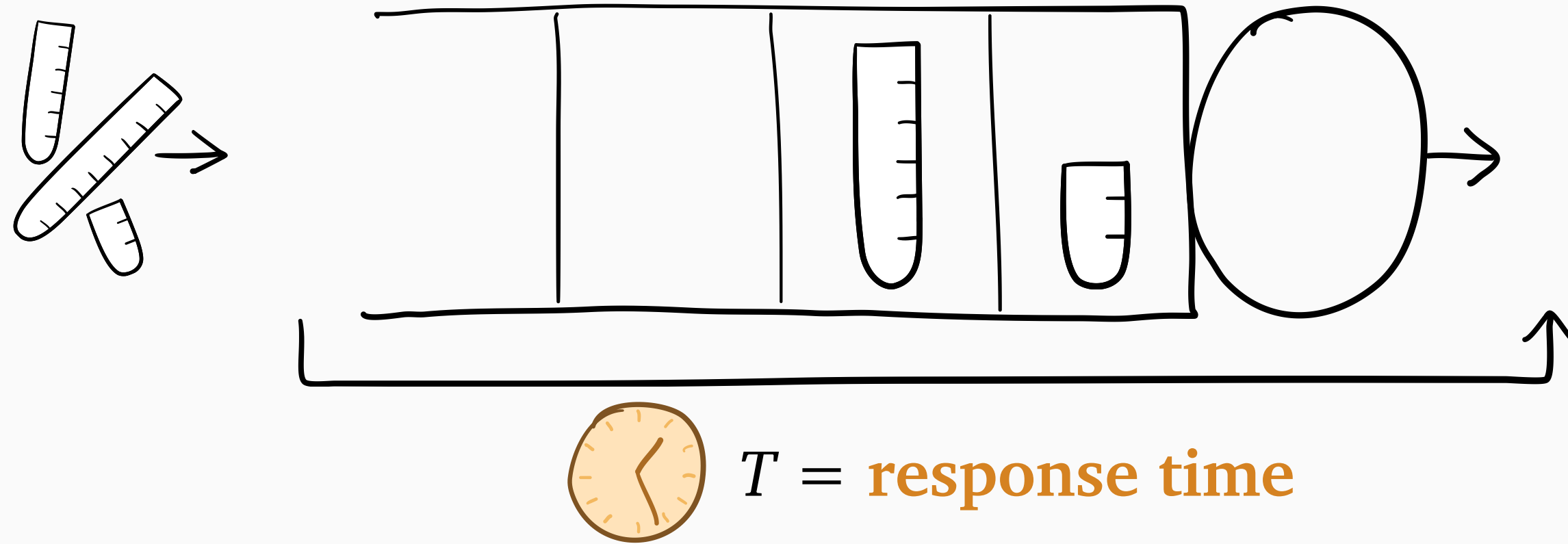


How should we schedule jobs to minimize delay?

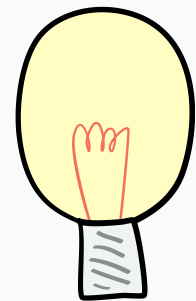


? Minimize $E[T]$?

How should we schedule jobs to minimize delay?

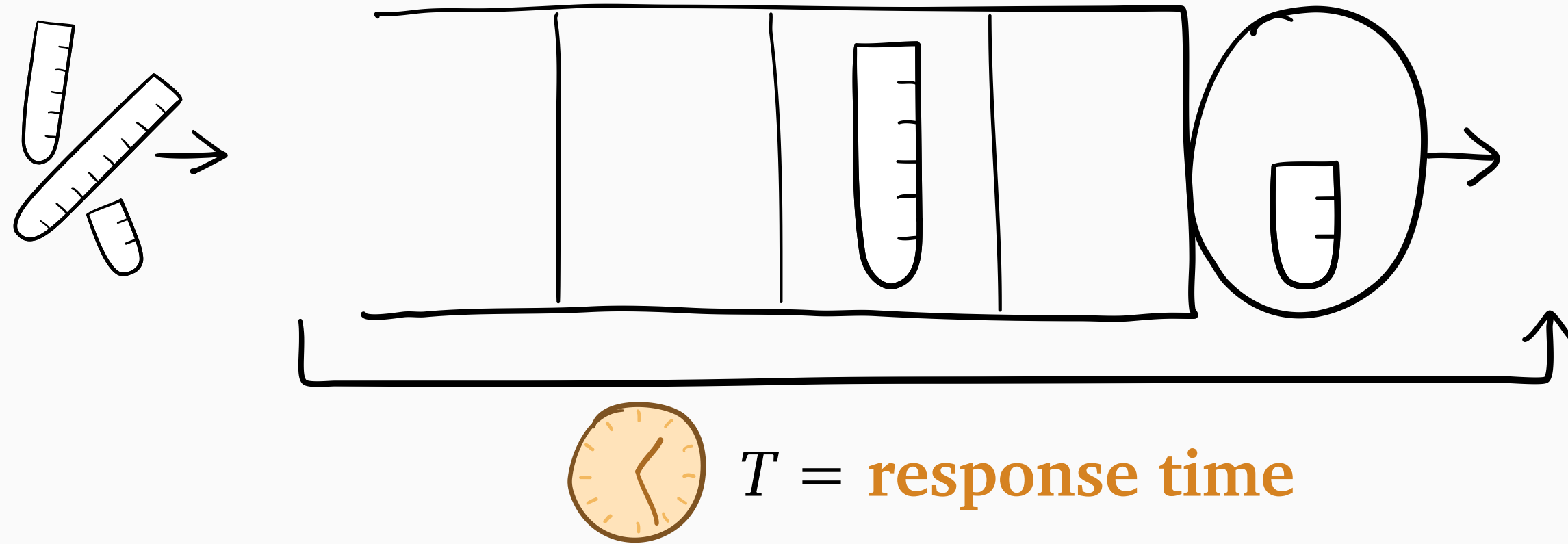


Minimize $E[T]$?



Serve short jobs
before long jobs

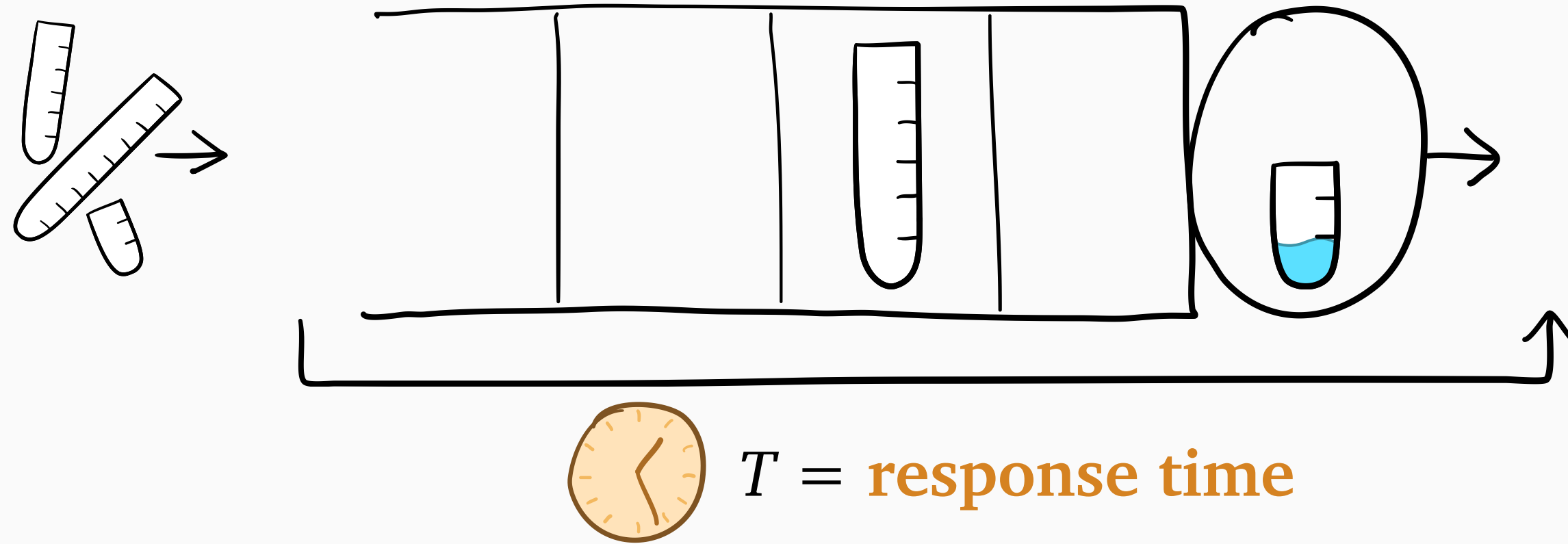
How should we schedule jobs to minimize delay?



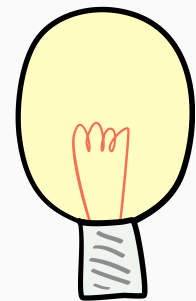
? Minimize $E[T]$?

💡 Serve short jobs
before long jobs

How should we schedule jobs to minimize delay?

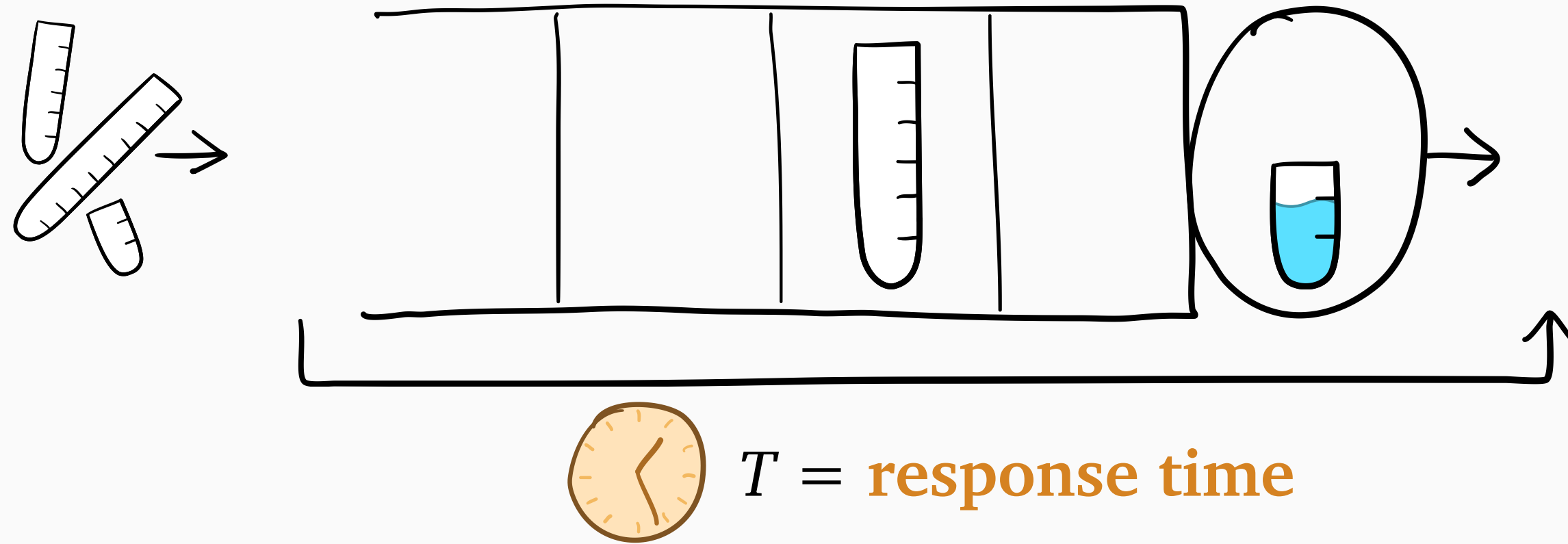


Minimize $E[T]$?



Serve short jobs
before long jobs

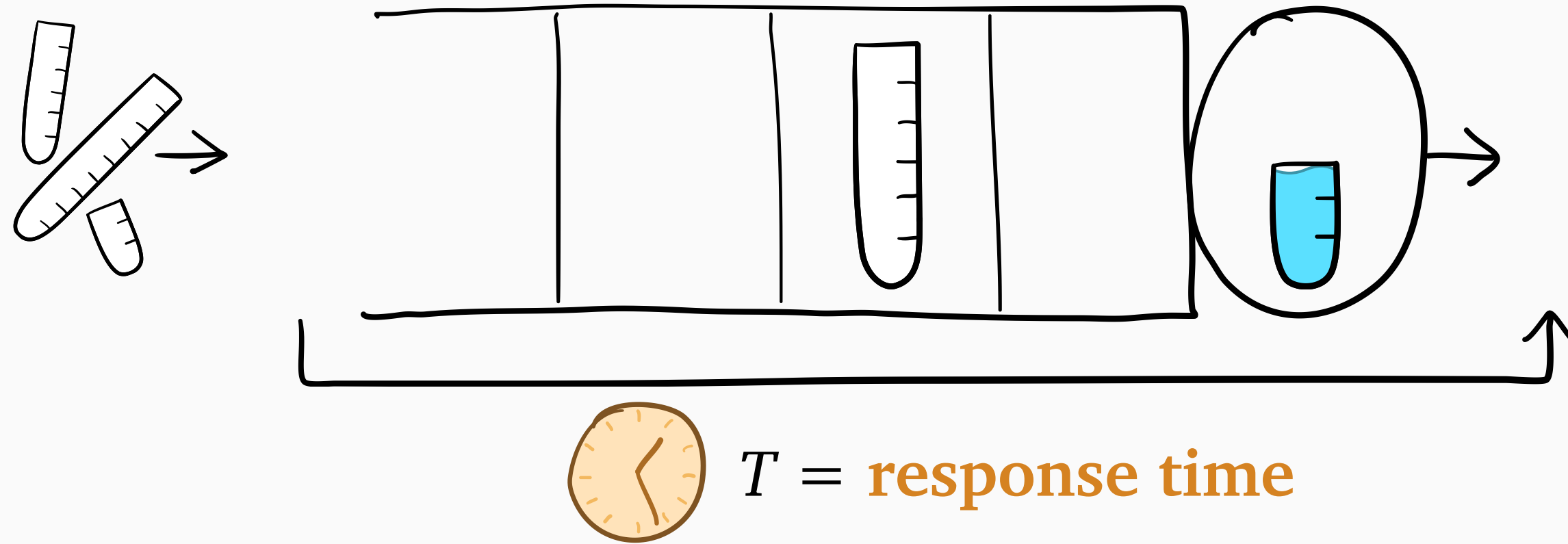
How should we schedule jobs to minimize delay?



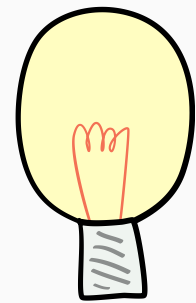
? Minimize $E[T]$?

💡 Serve short jobs
before long jobs

How should we schedule jobs to minimize delay?

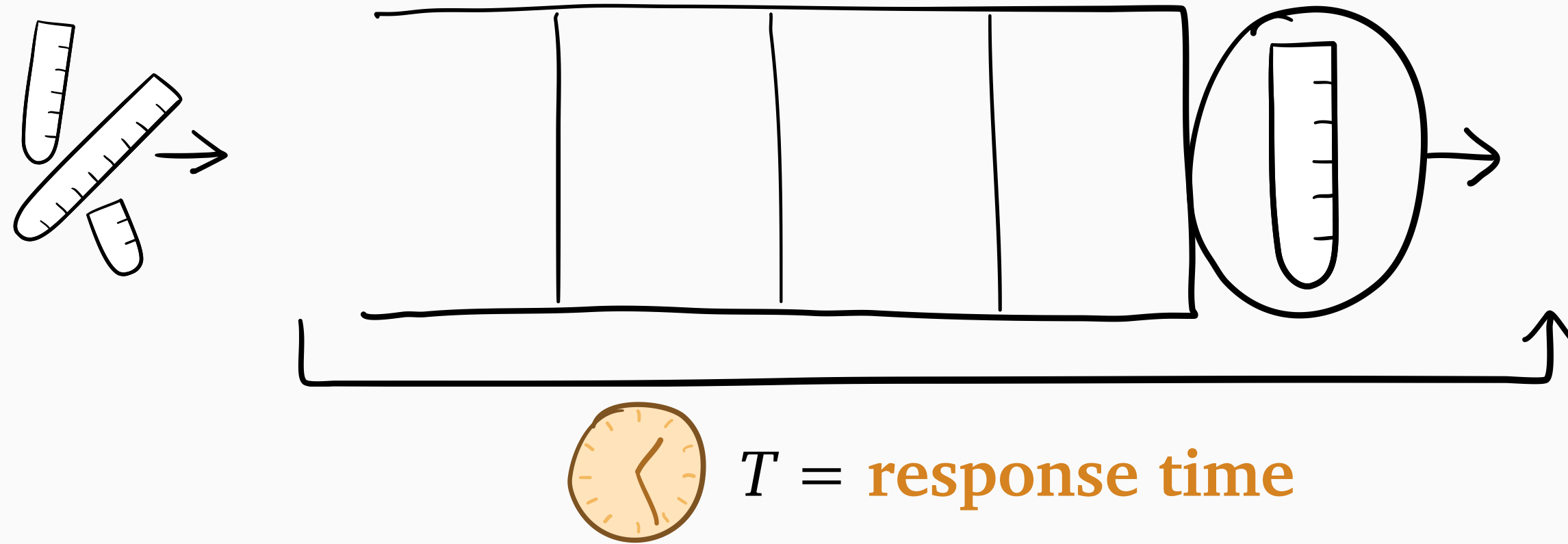


Minimize $E[T]$?

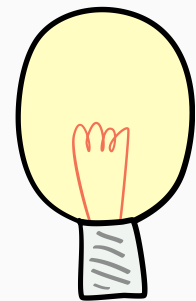


Serve short jobs
before long jobs

How should we schedule jobs to minimize delay?

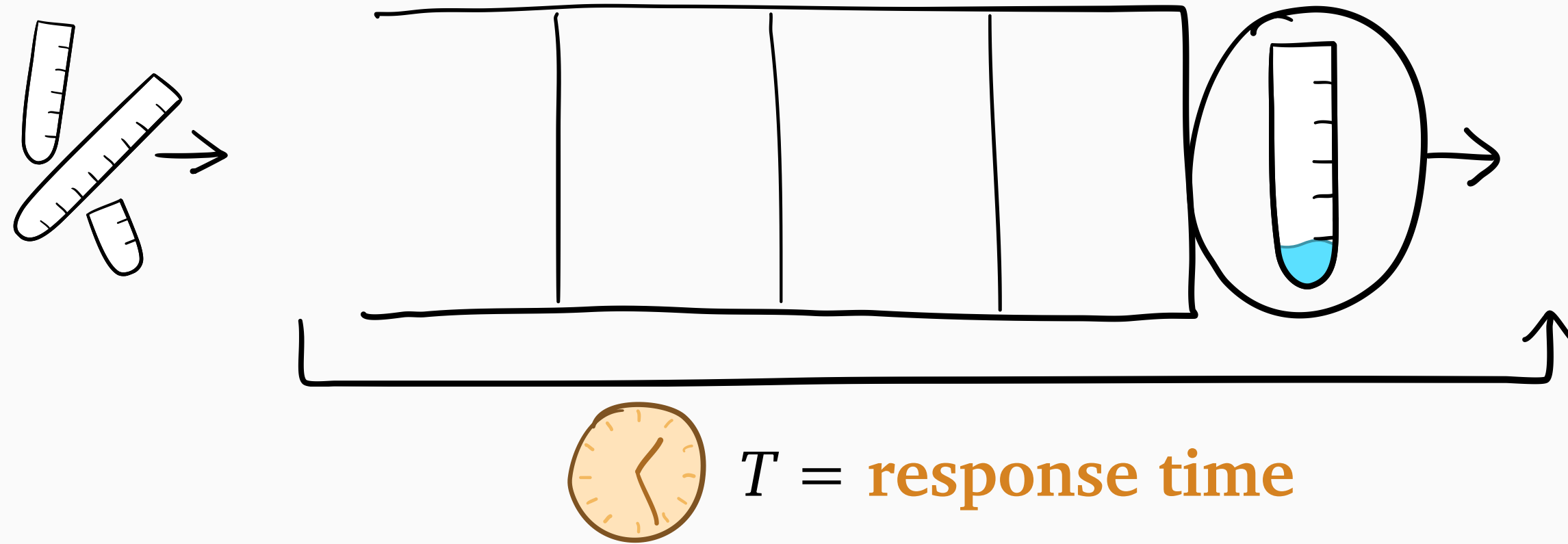


Minimize $E[T]$?



Serve short jobs
before long jobs

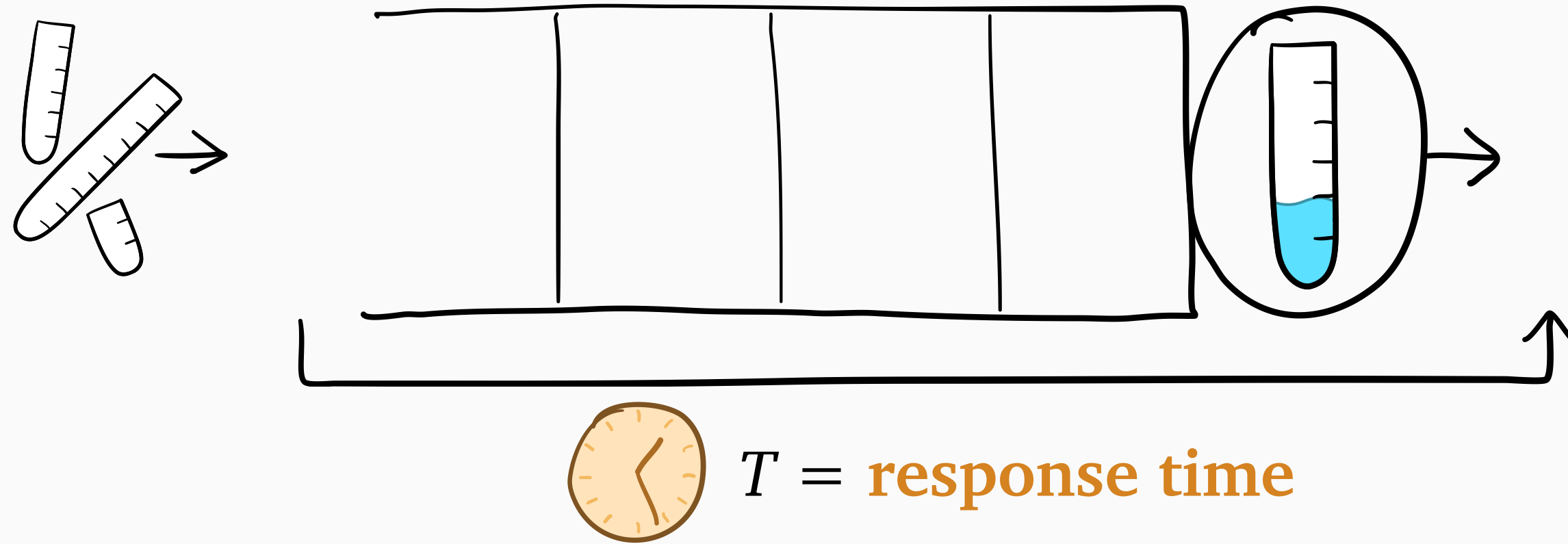
How should we schedule jobs to minimize delay?



? Minimize $E[T]$?

💡 Serve short jobs
before long jobs

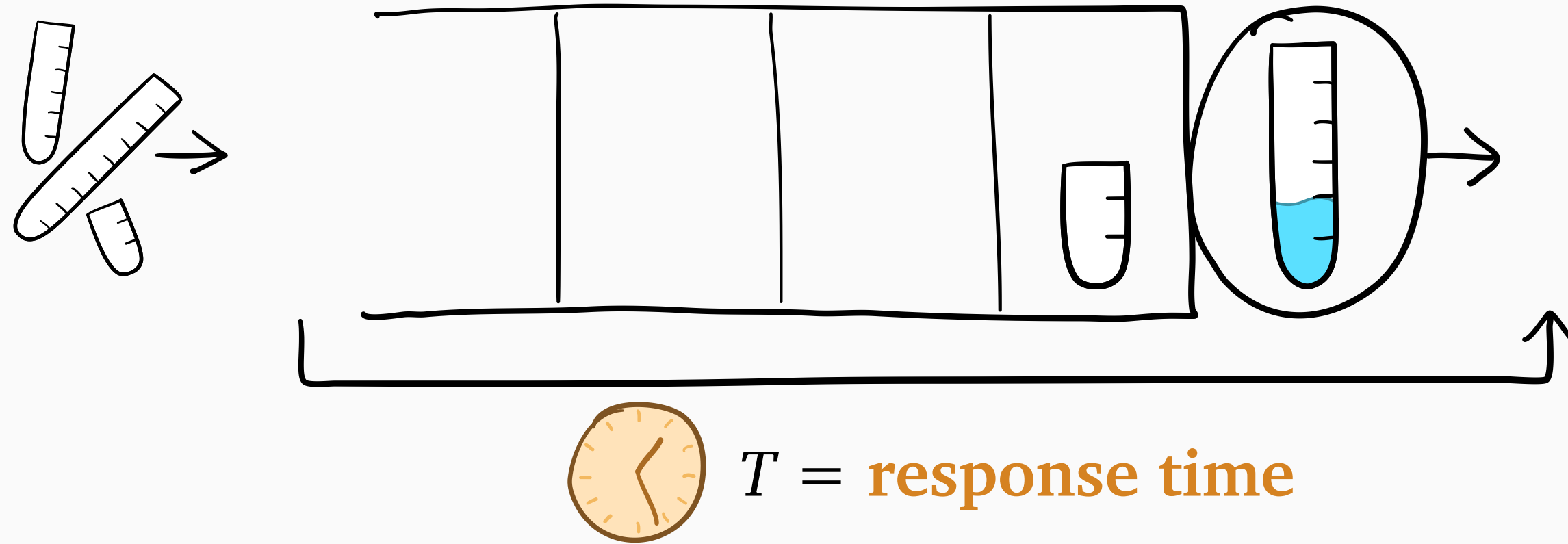
How should we schedule jobs to minimize delay?



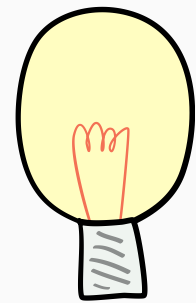
? Minimize $E[T]$?

💡 Serve short jobs
before long jobs

How should we schedule jobs to minimize delay?

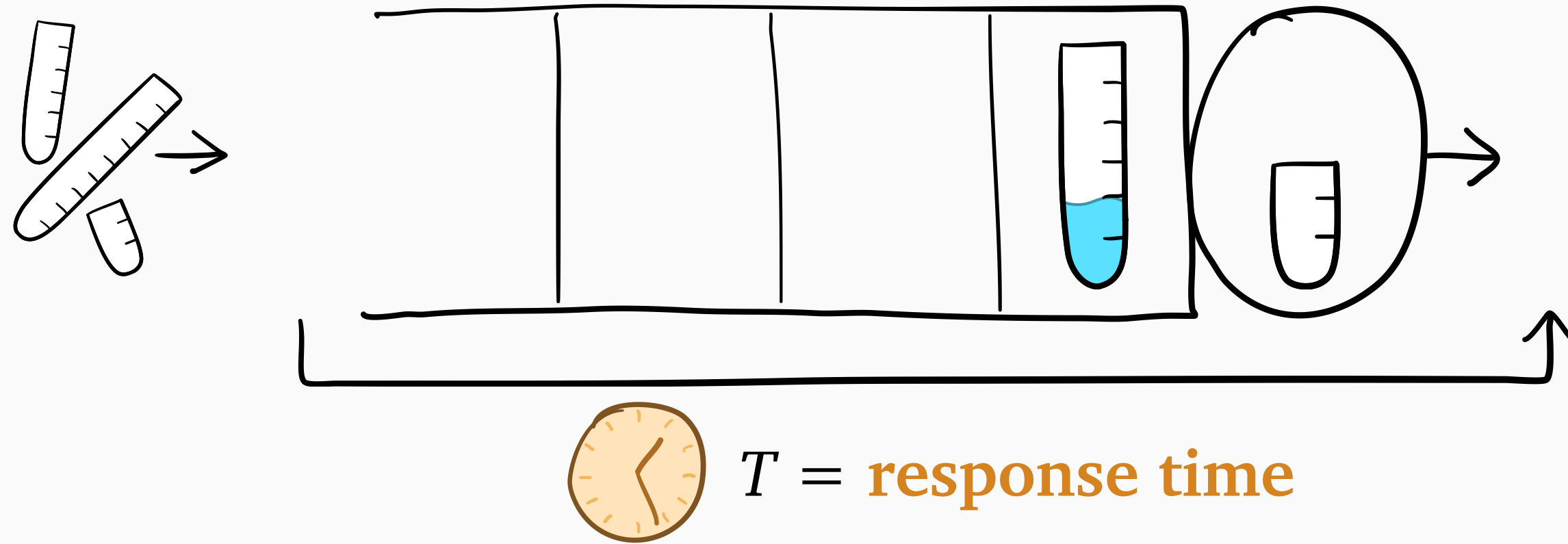


Minimize $E[T]$?



Serve short jobs
before long jobs

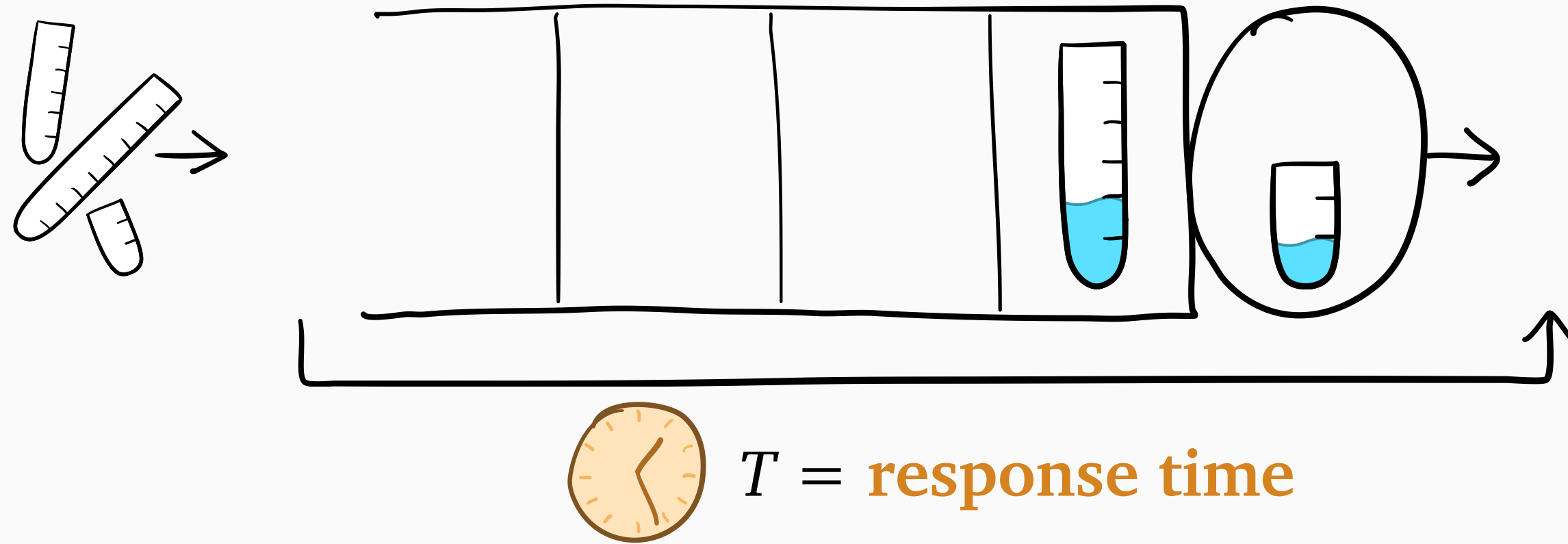
How should we schedule jobs to minimize delay?



? Minimize $E[T]$?

💡 Serve short jobs
before long jobs

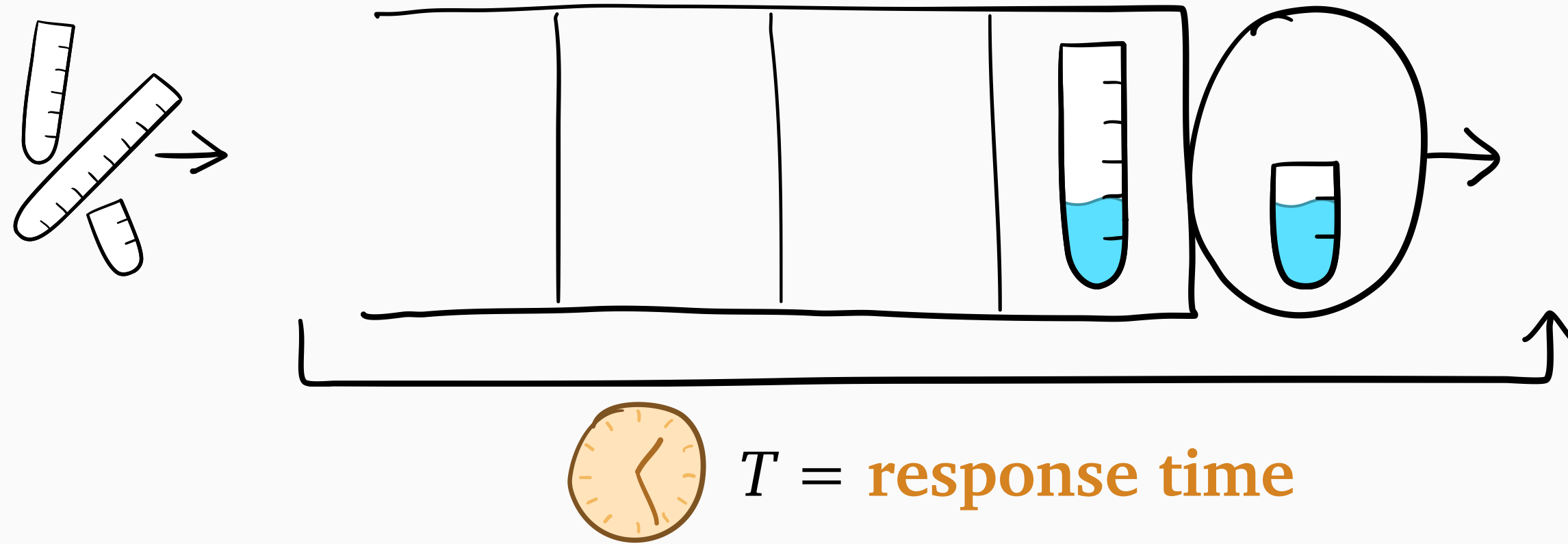
How should we schedule jobs to minimize delay?



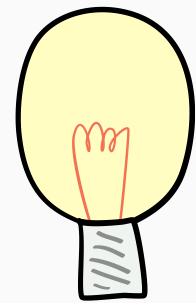
? Minimize $E[T]$?

💡 Serve short jobs
before long jobs

How should we schedule jobs to minimize delay?

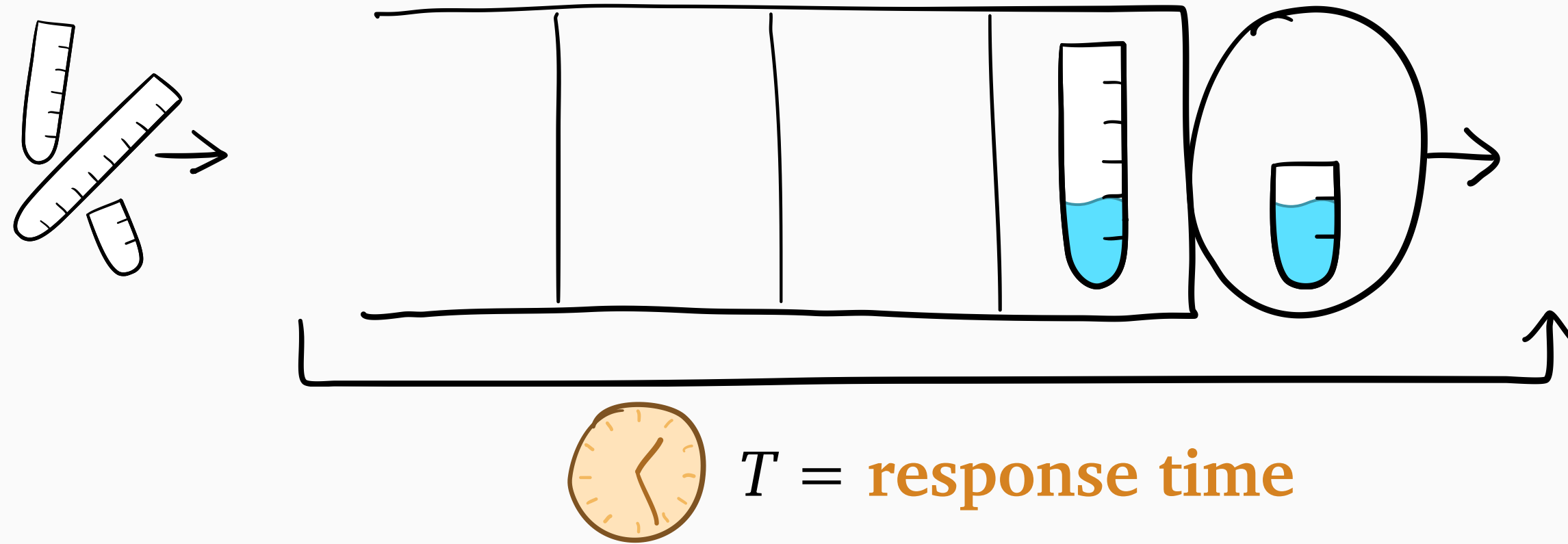


Minimize $E[T]$?



Serve short jobs
before long jobs

How should we schedule jobs to minimize delay?

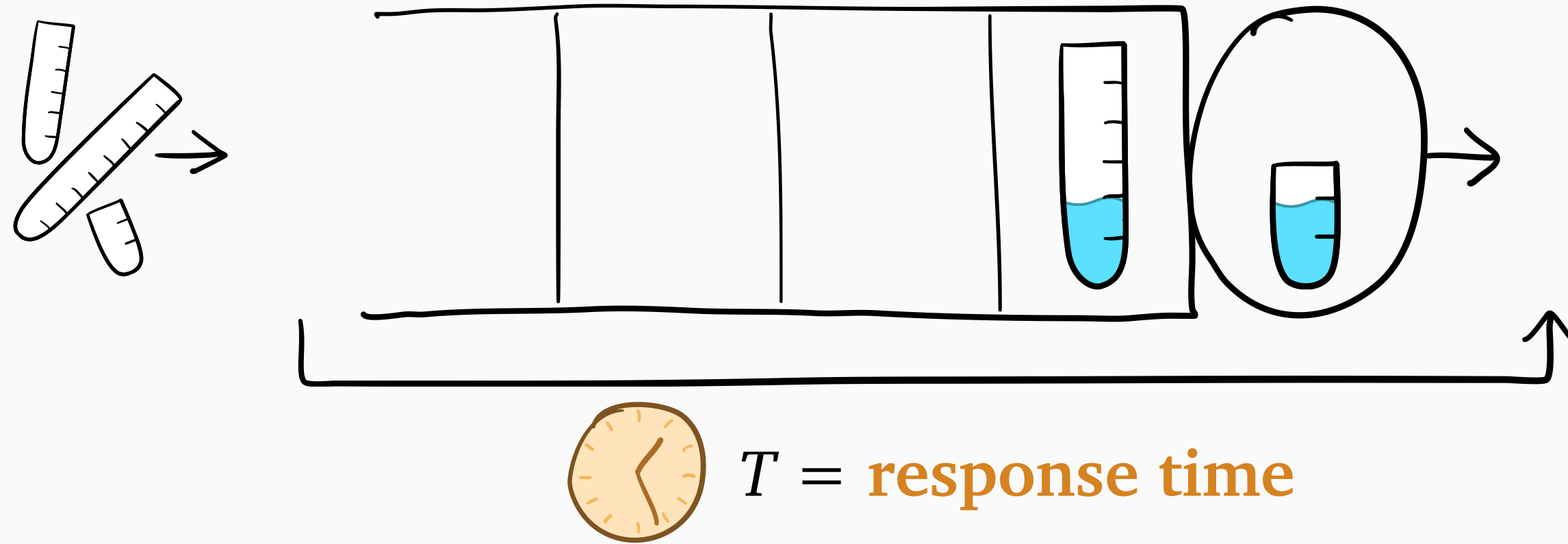


? Minimize $E[T]$?

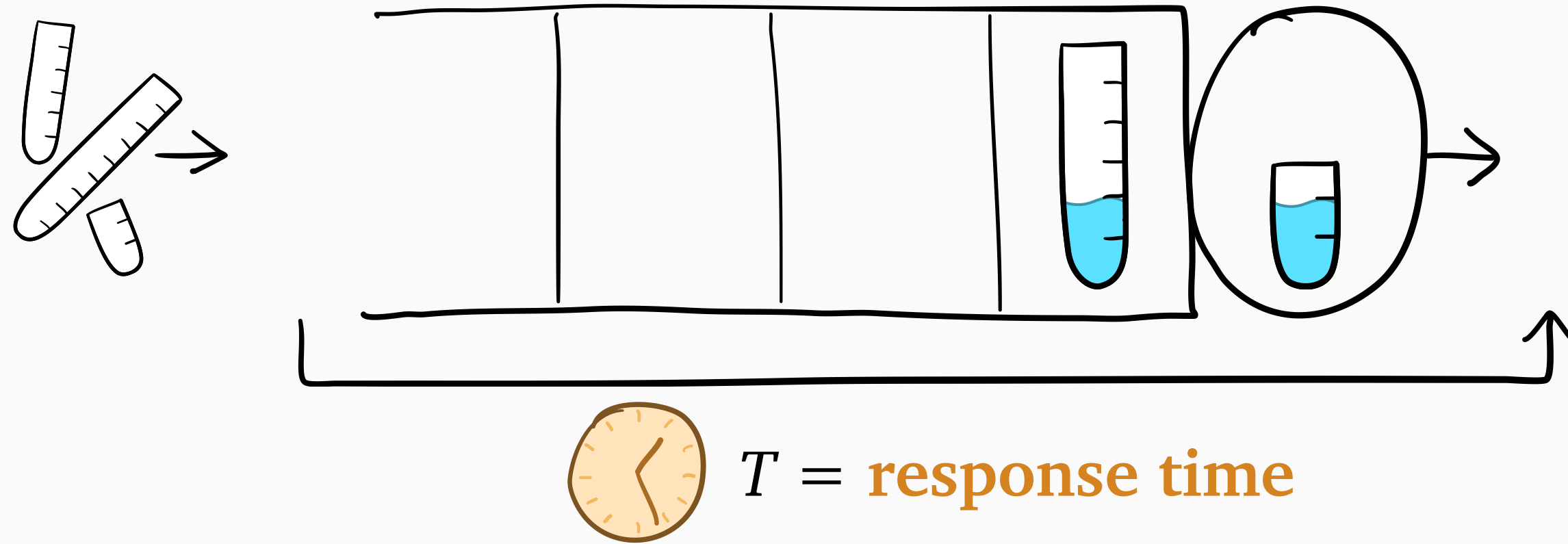
💡 Serve short jobs
before long jobs

🏆 **SRPT**: minimizes $E[T]$
shortest remaining
processing time

TCS vs. Queueing



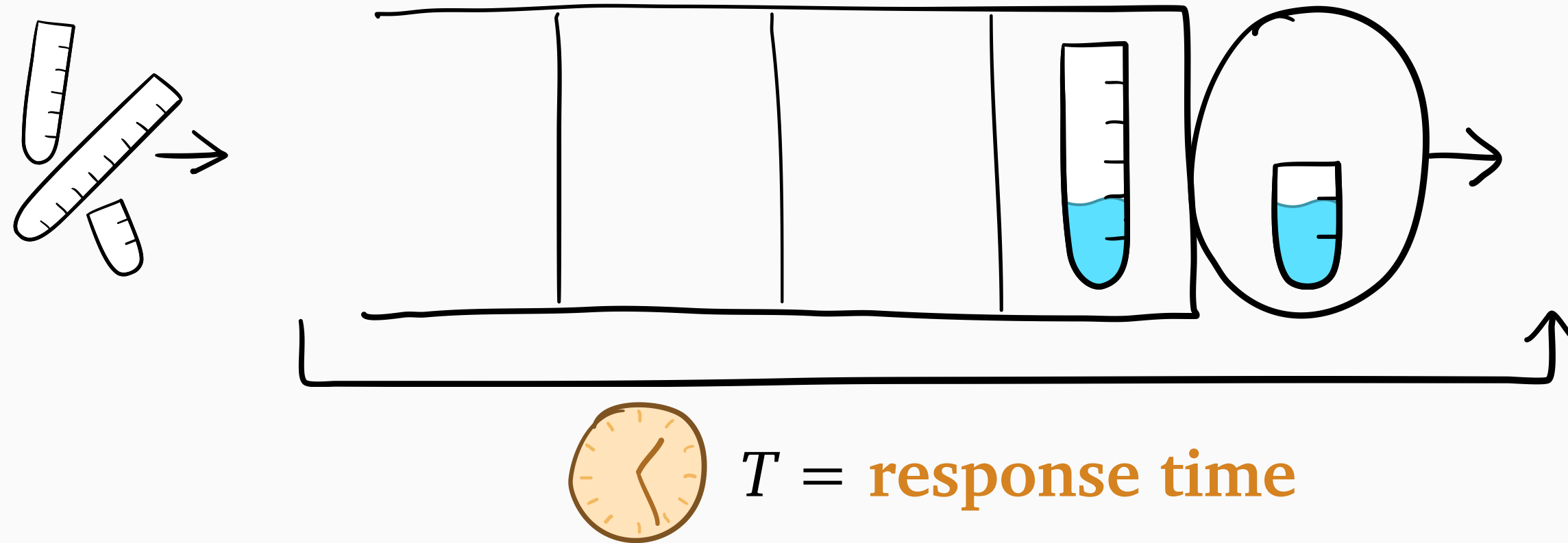
TCS vs. Queueing



TCS

n arbitrary arrivals
 T is tuple of n times

TCS vs. Queueing



TCS

n arbitrary arrivals
 T is tuple of n times

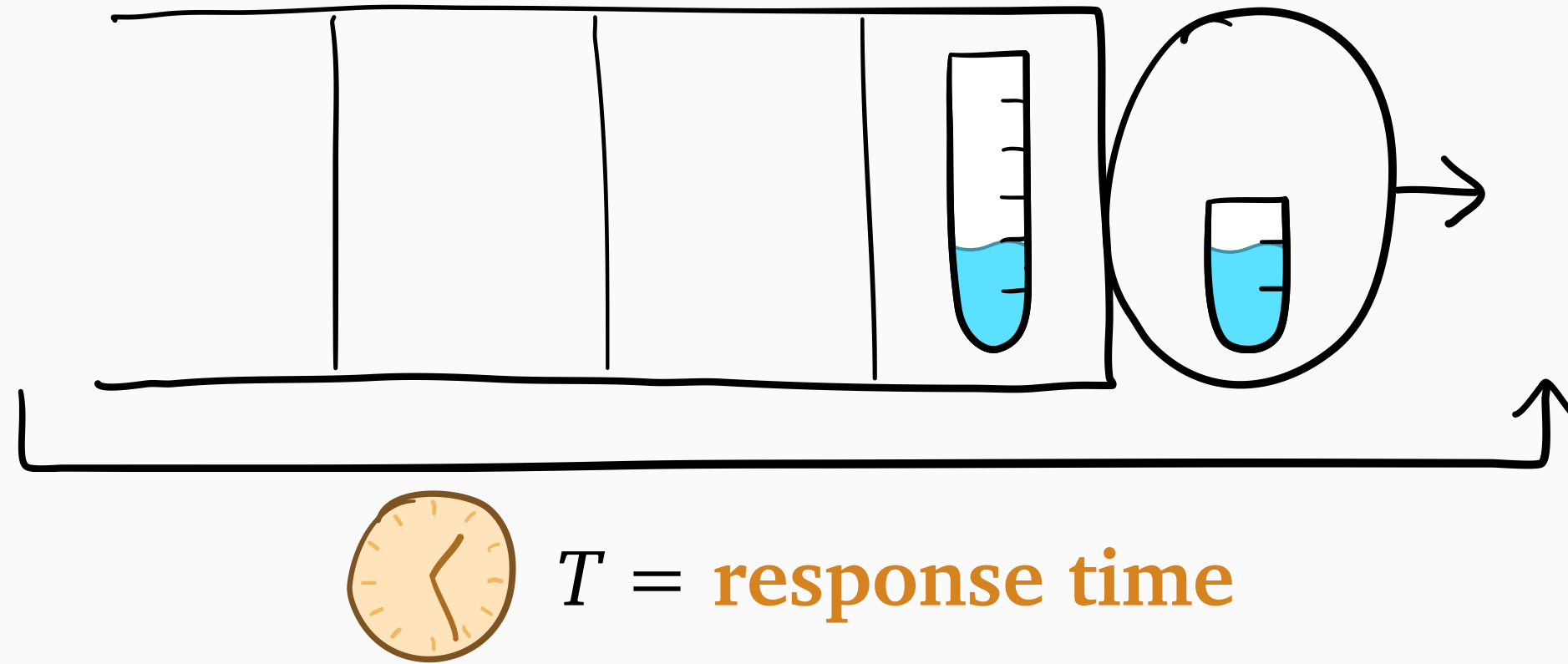
Queueing

infinite stochastic sequence of arrivals
 T is a limiting distribution

TCS vs. Queueing

M/G arrivals

- arrival rate λ (Poisson)
- job size dist. S



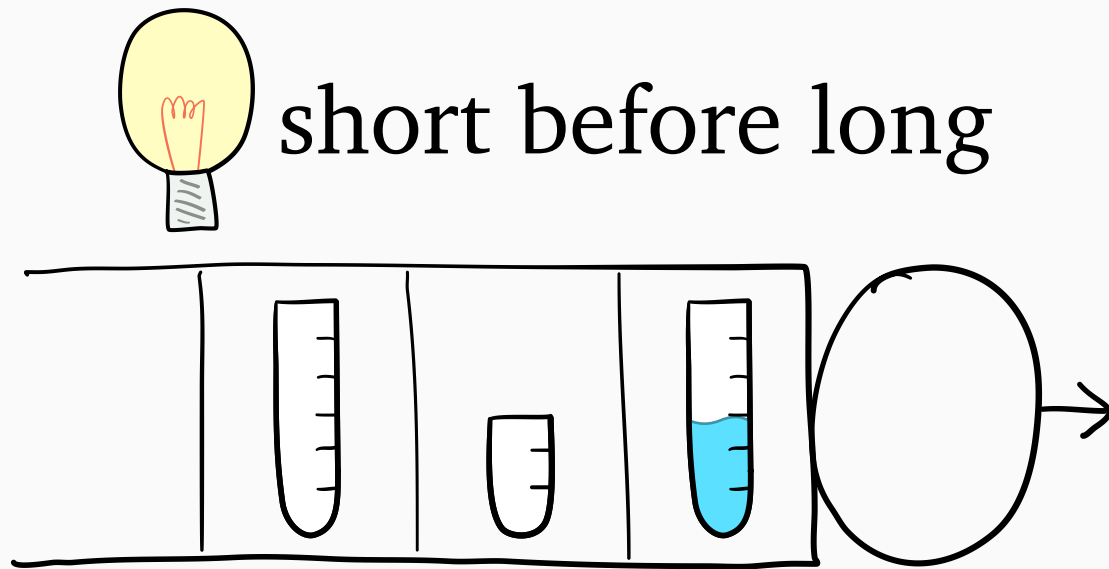
TCS

n arbitrary arrivals
 T is tuple of n times

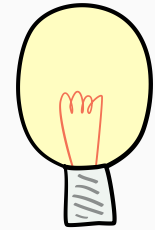
Queueing

infinite stochastic sequence of arrivals
 T is a limiting distribution

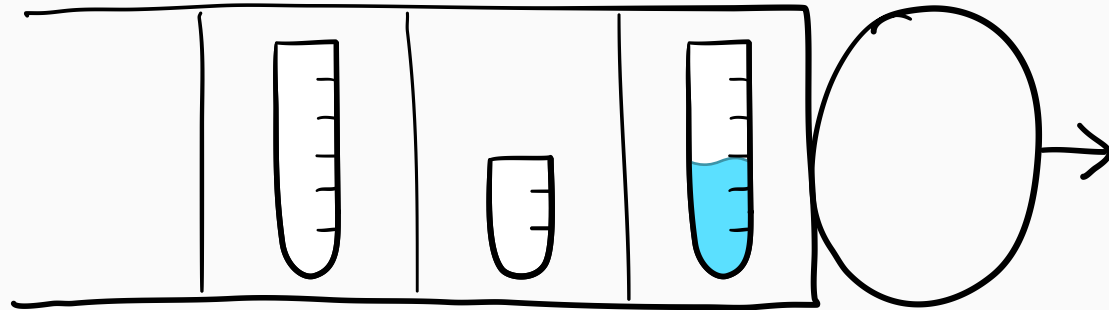
Job size uncertainty



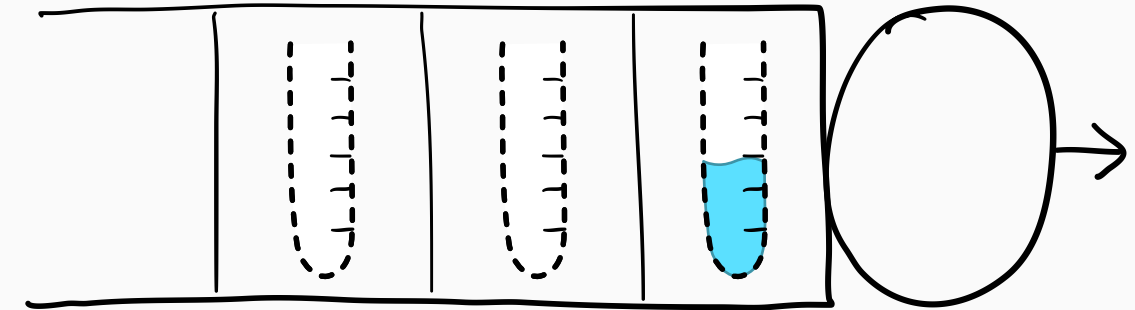
Job size uncertainty



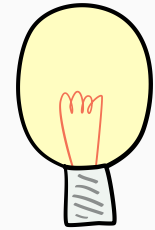
short before long



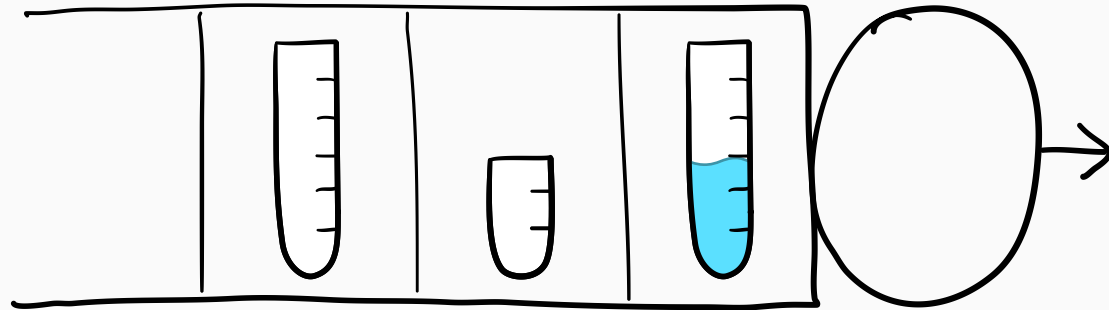
unknown sizes



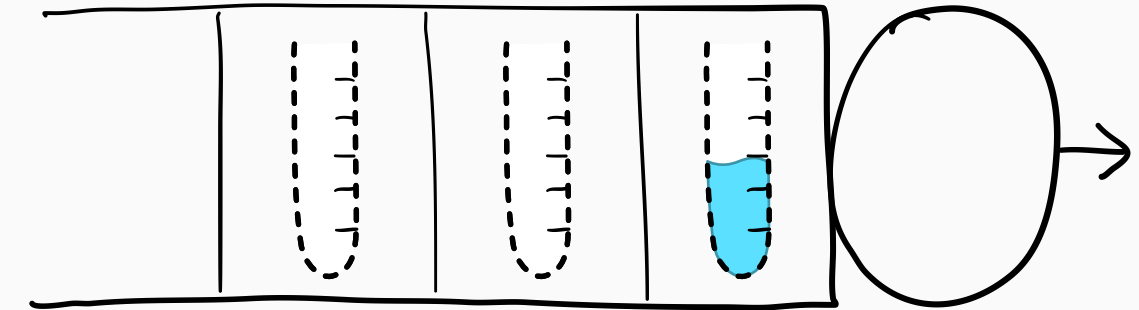
Job size uncertainty



short before long

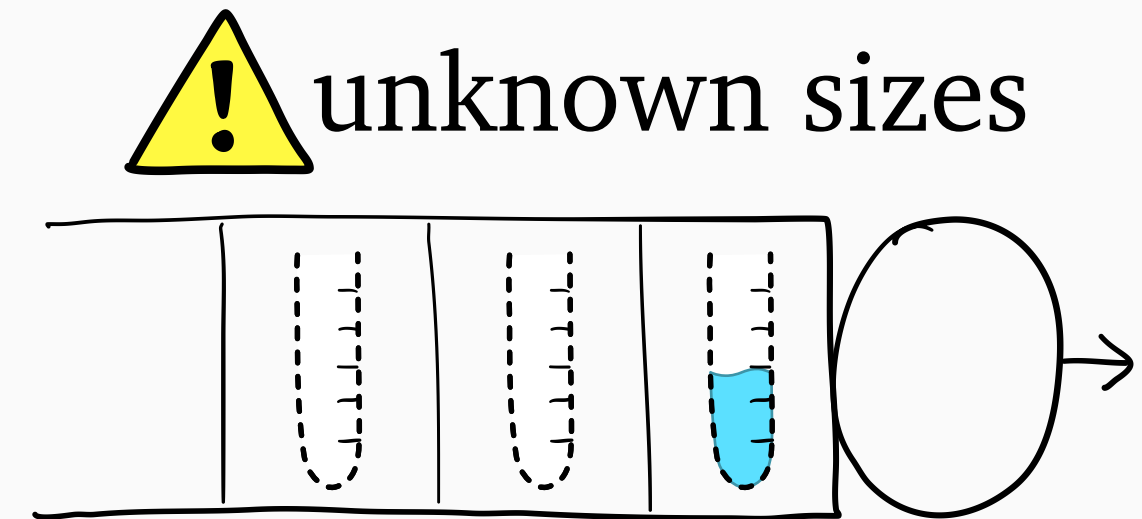
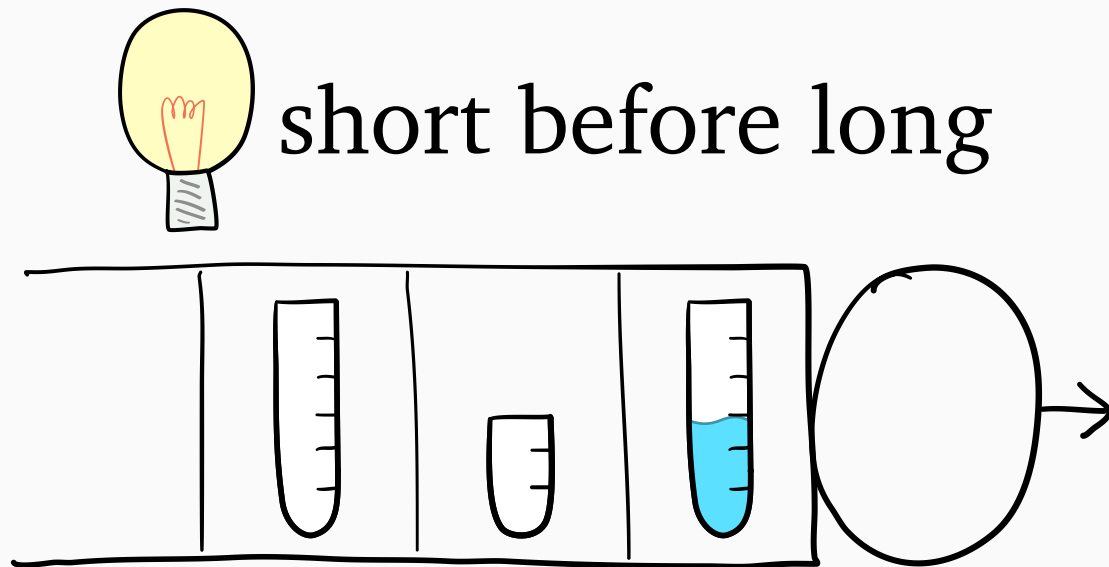


unknown sizes



What info can we use?

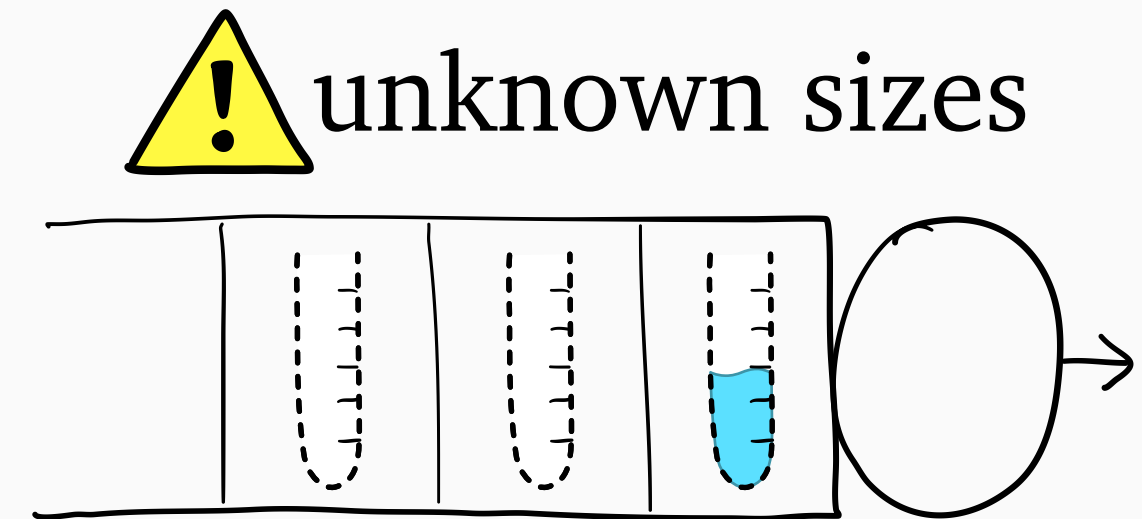
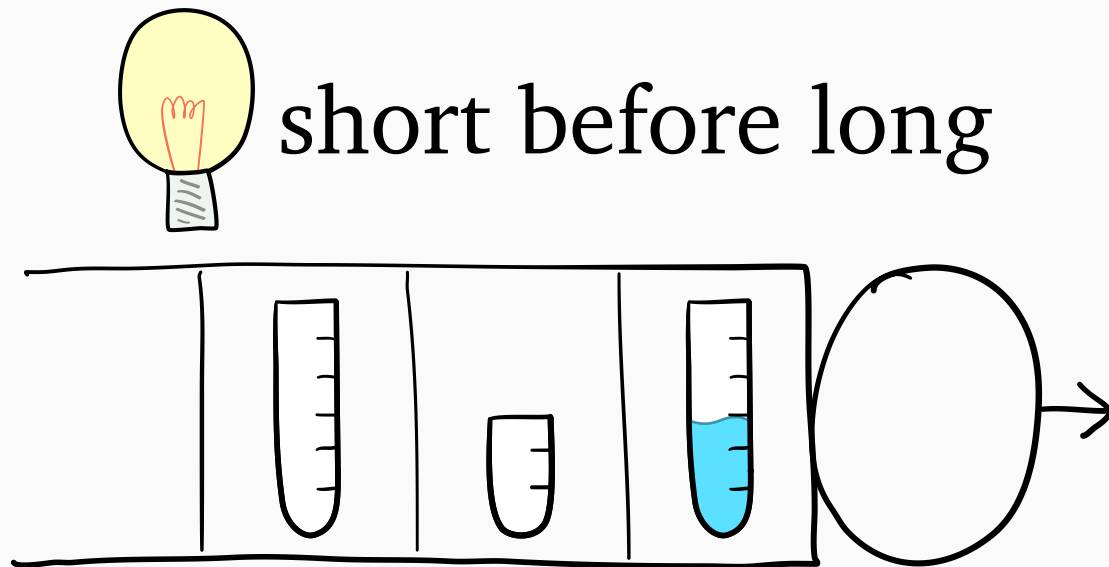
Job size uncertainty



? What info can we use?

- job **ages**

Job size uncertainty



- ?
- What info can we use?
- job **ages**
 - size distribution S

Scheduling with unknown job sizes



What info can we use?

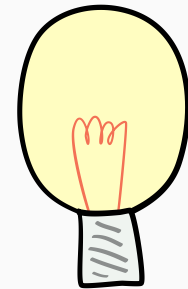
- job **ages**
- size distribution S

Scheduling with unknown job sizes



What info can we use?

- job **ages**
- size distribution S



Gittins policy construction:

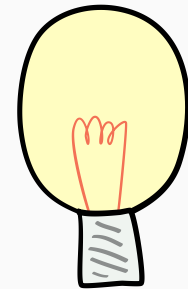
size distribution $S \mapsto$ policy **Gittins**(S)

Scheduling with unknown job sizes



What info can we use?

- job **ages**
- size distribution S



Gittins policy construction:

size distribution $S \mapsto$ policy **Gittins**(S)

age $a \mapsto$ priority **rank**(a)

Outline of Part 1

? What is **Gittins**?

main focus

? Why are **Gittins** (and **SRPT**) optimal?

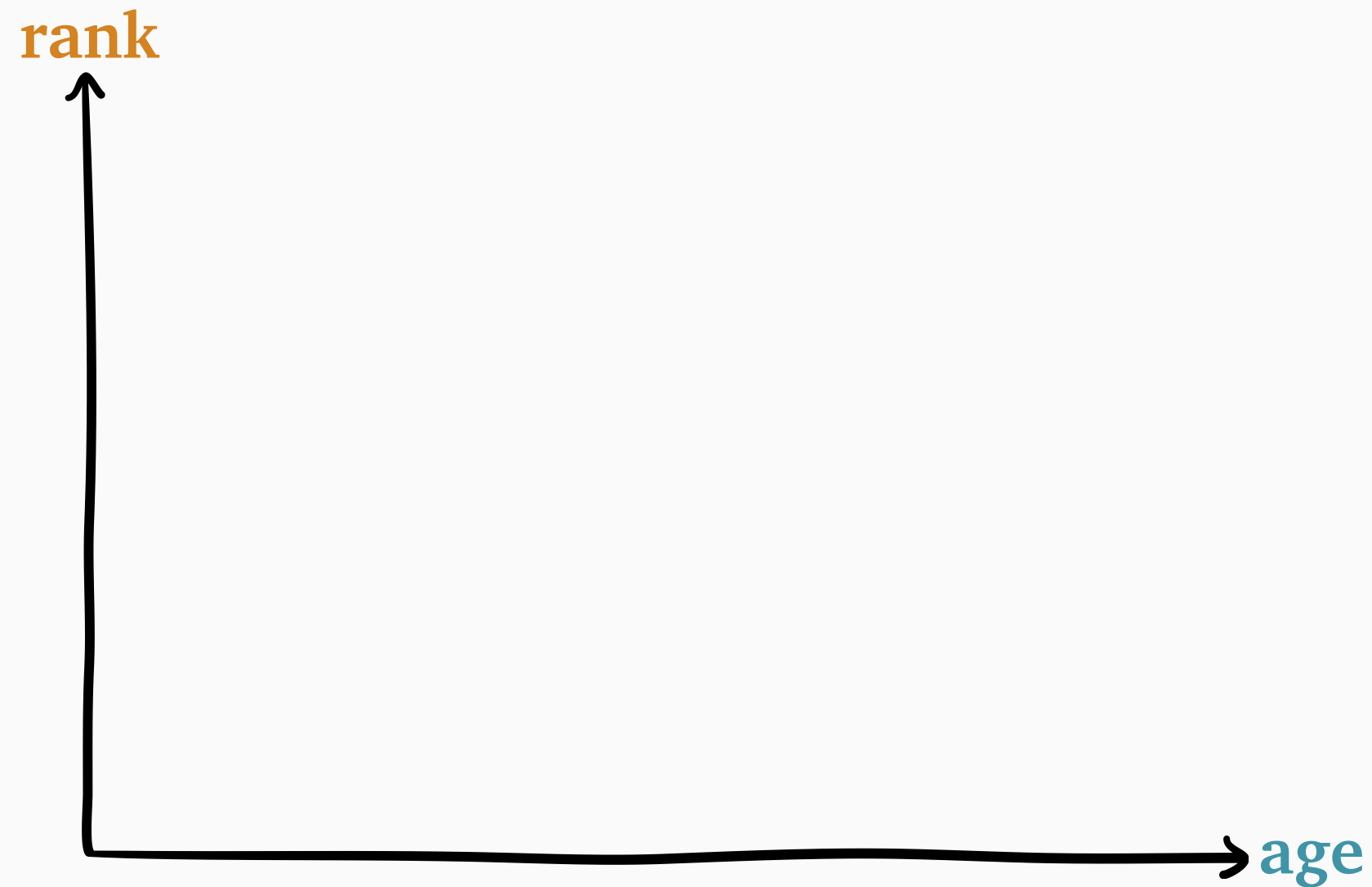
? **Predictions:** what if we don't know exact distributions (or sizes)?

Scheduling with **rank** functions

Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

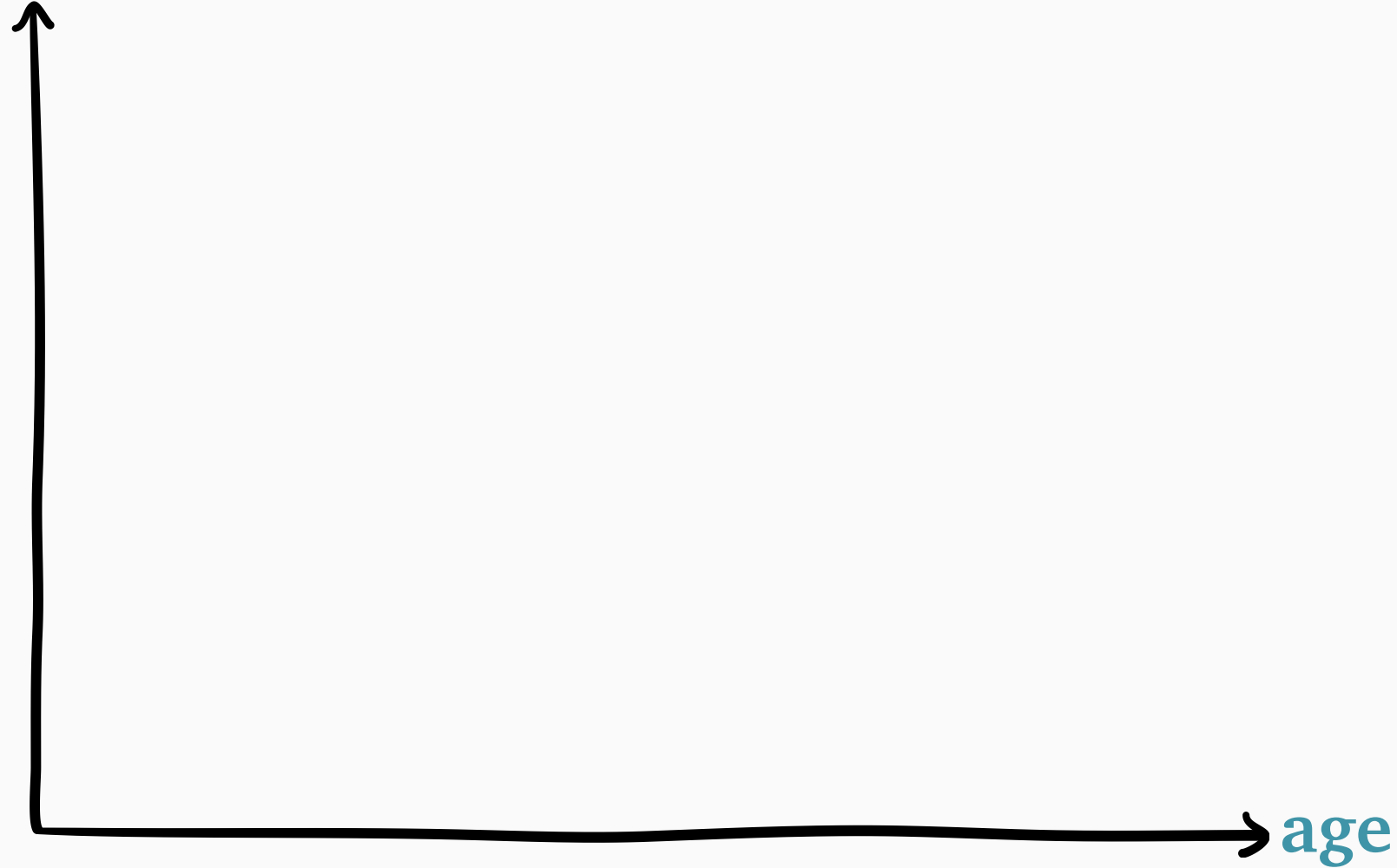


Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

$$\text{rank}(a) = E[S - a \mid S > a]$$

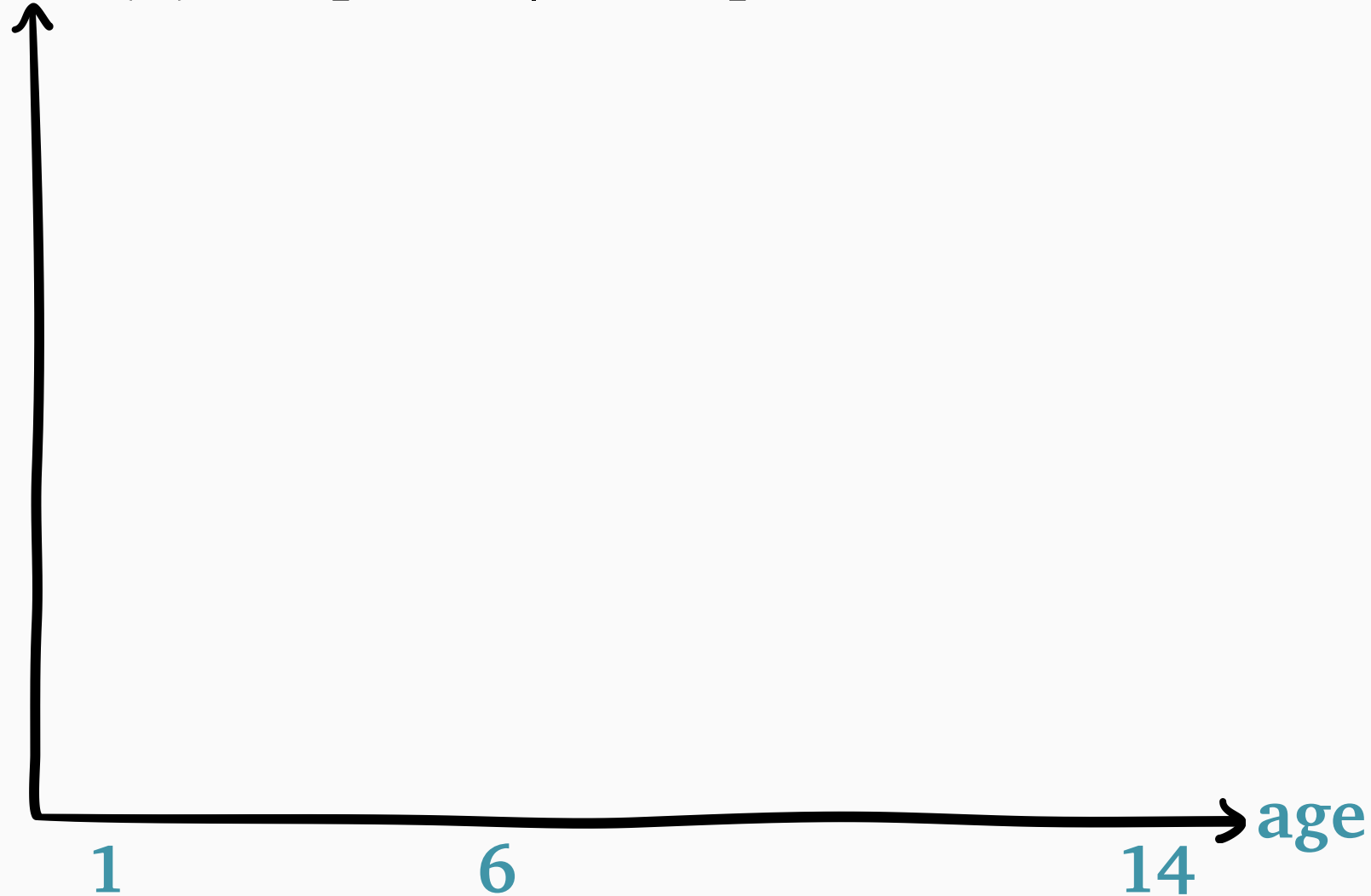


Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

$$\text{rank}(a) = E[S - a \mid S > a]$$



Job size distribution:

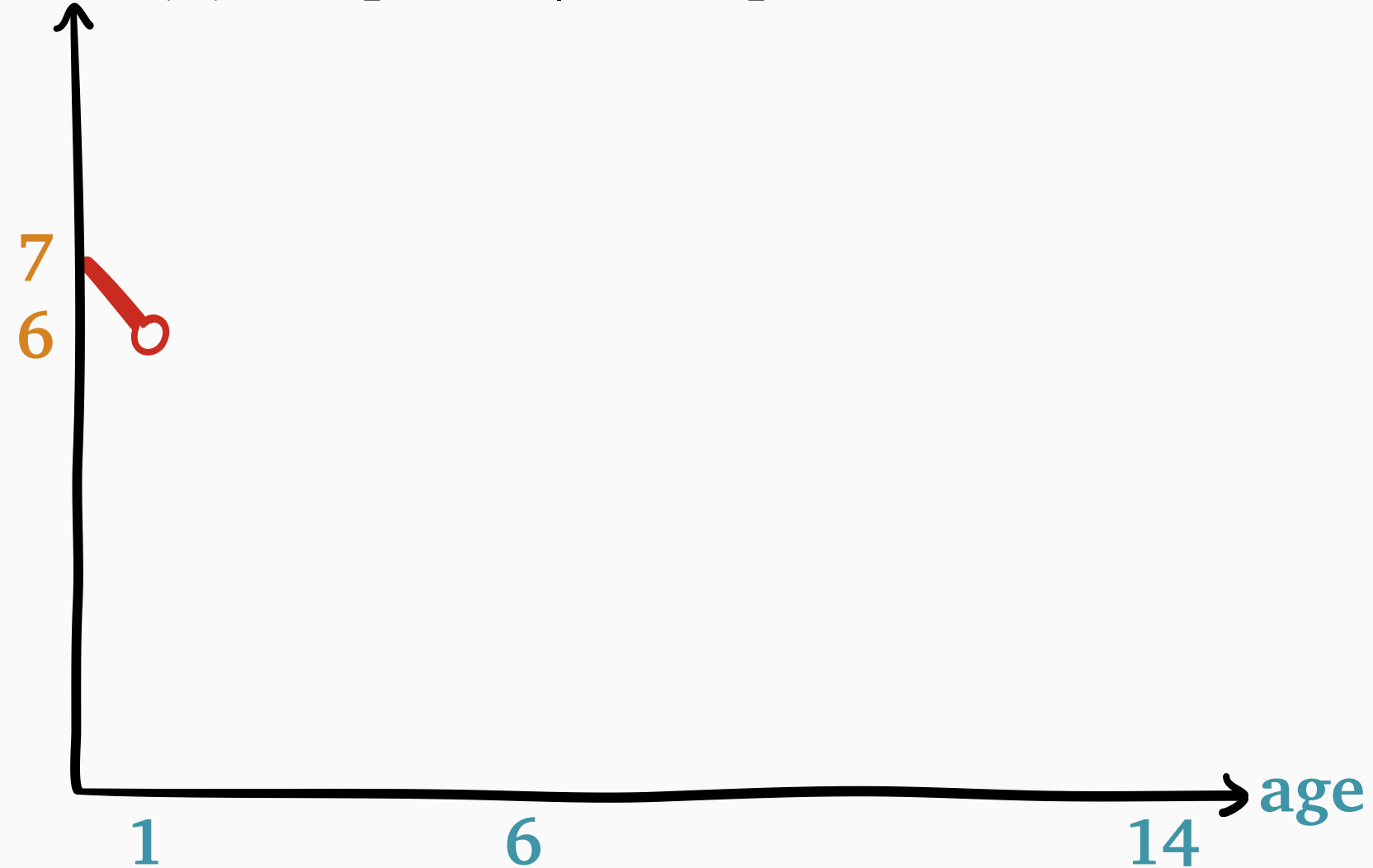
$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$

Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

$$\text{rank}(a) = E[S - a \mid S > a]$$



Job size distribution:

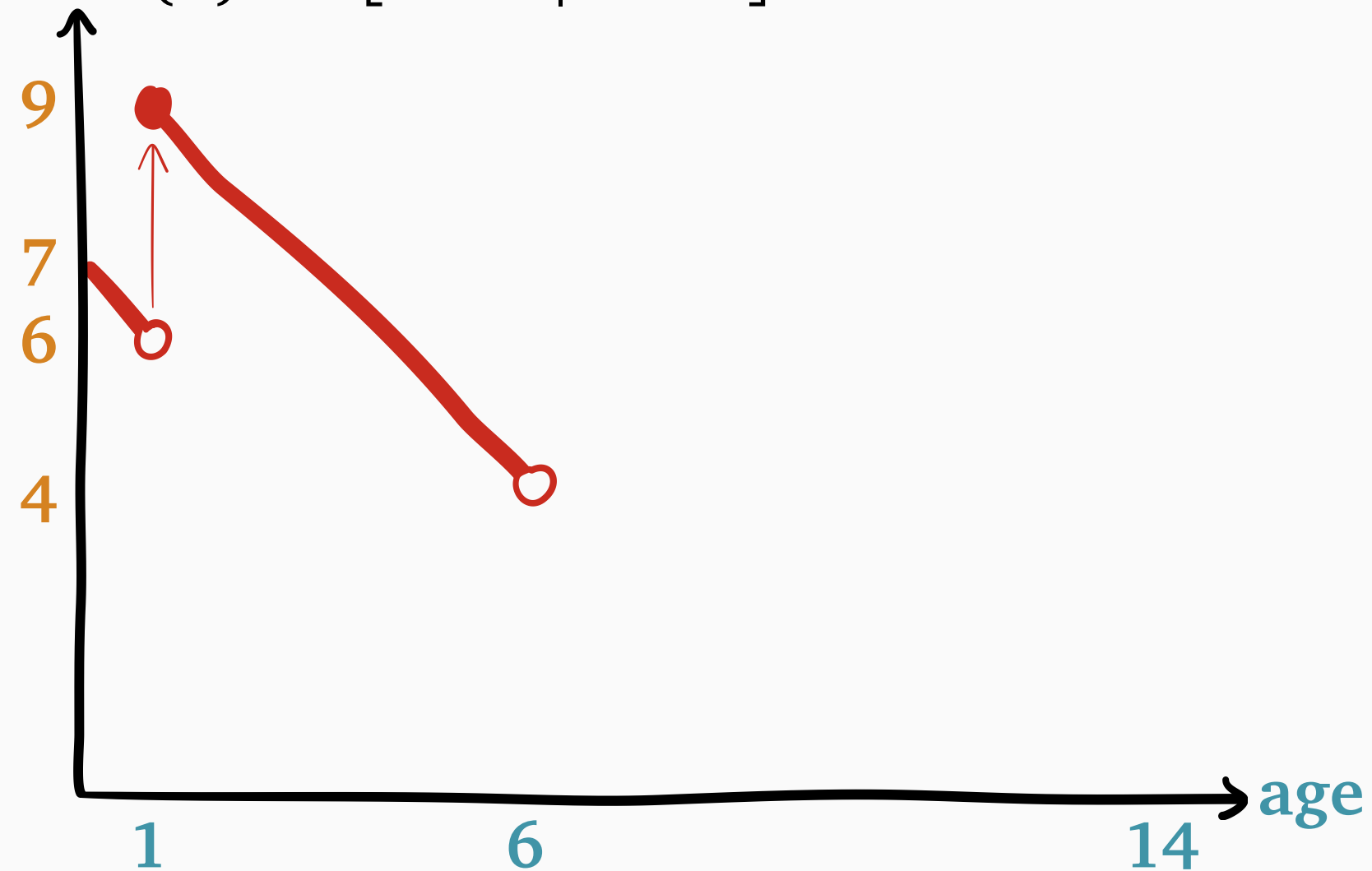
$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$

Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

$$\text{rank}(a) = E[S - a \mid S > a]$$



Job size distribution:

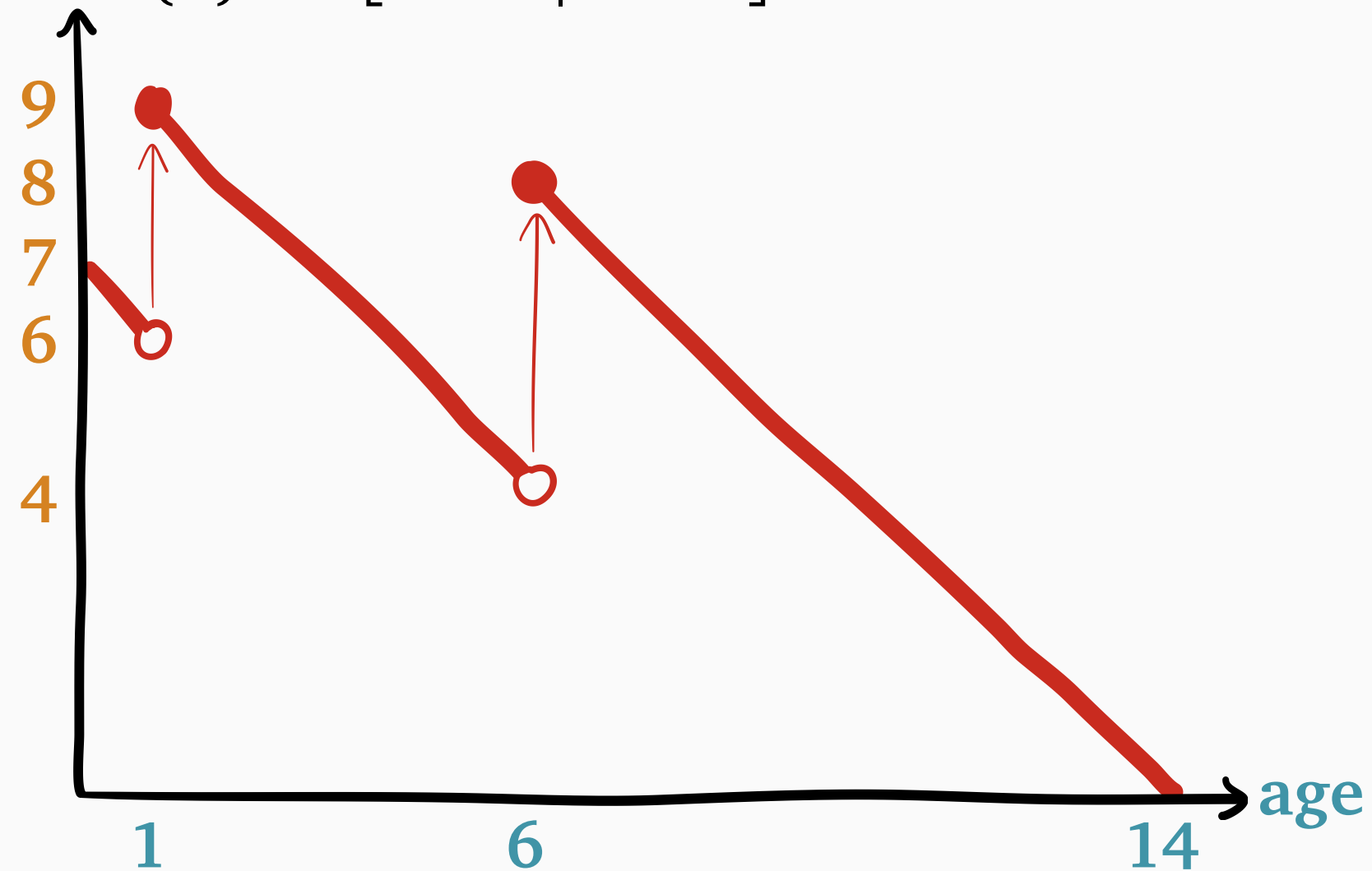
$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$

Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

$$\text{rank}(a) = E[S - a \mid S > a]$$



Job size distribution:

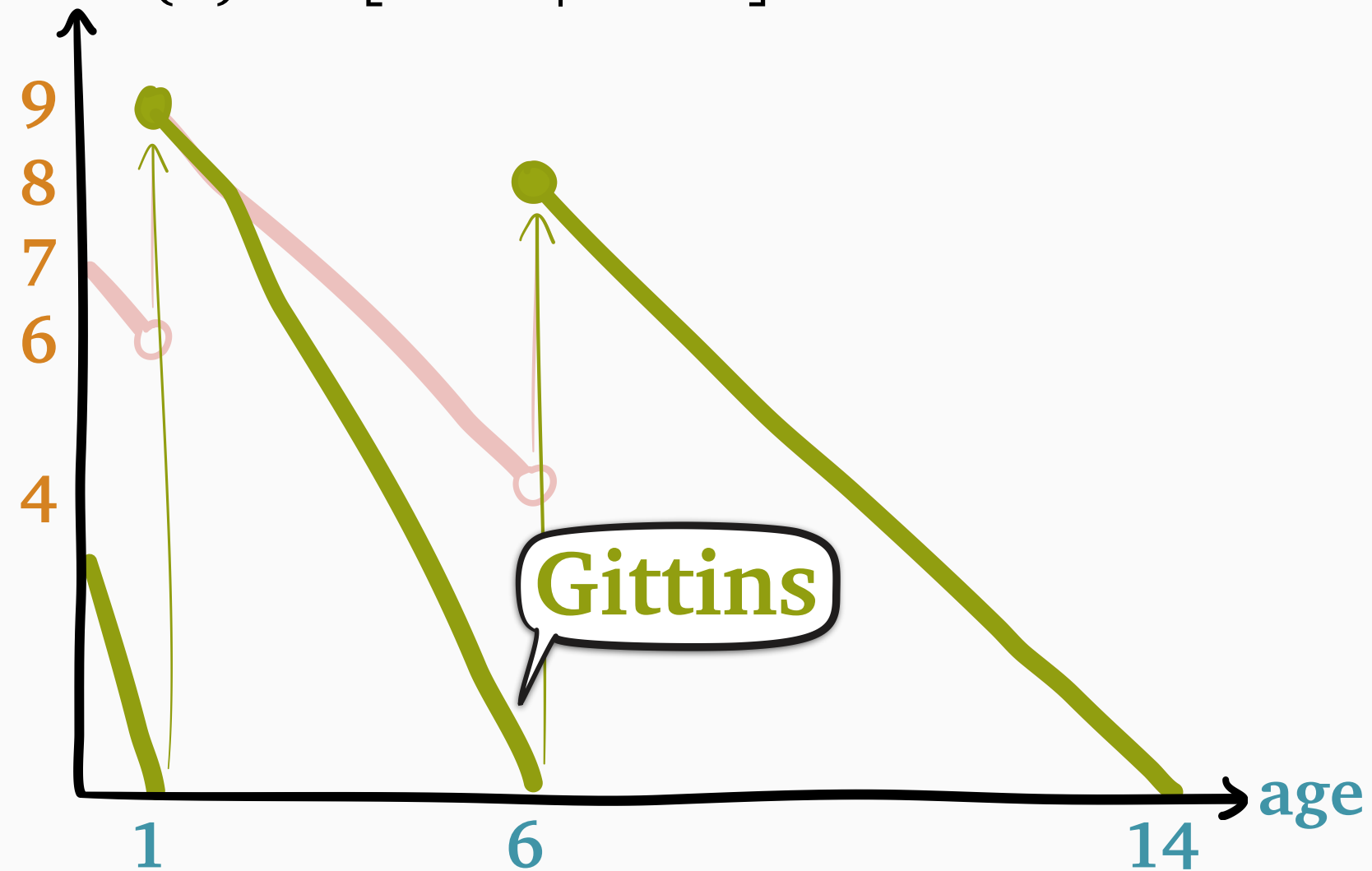
$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$

Scheduling with **rank** functions

SERPT

shortest *expected* remaining processing time

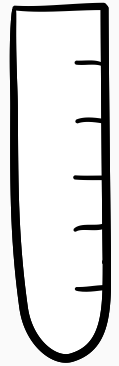
$$\text{rank}(a) = E[S - a \mid S > a]$$



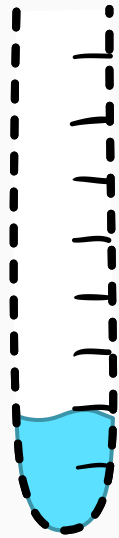
Job size distribution:

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

Defining the Gittins rank

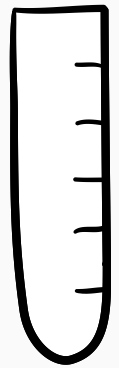


known
size r

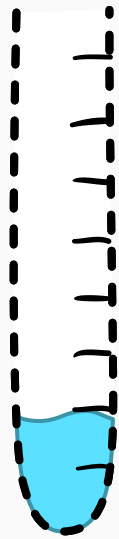


unknown size
($S \mid S > \text{age}$)

Defining the Gittins rank



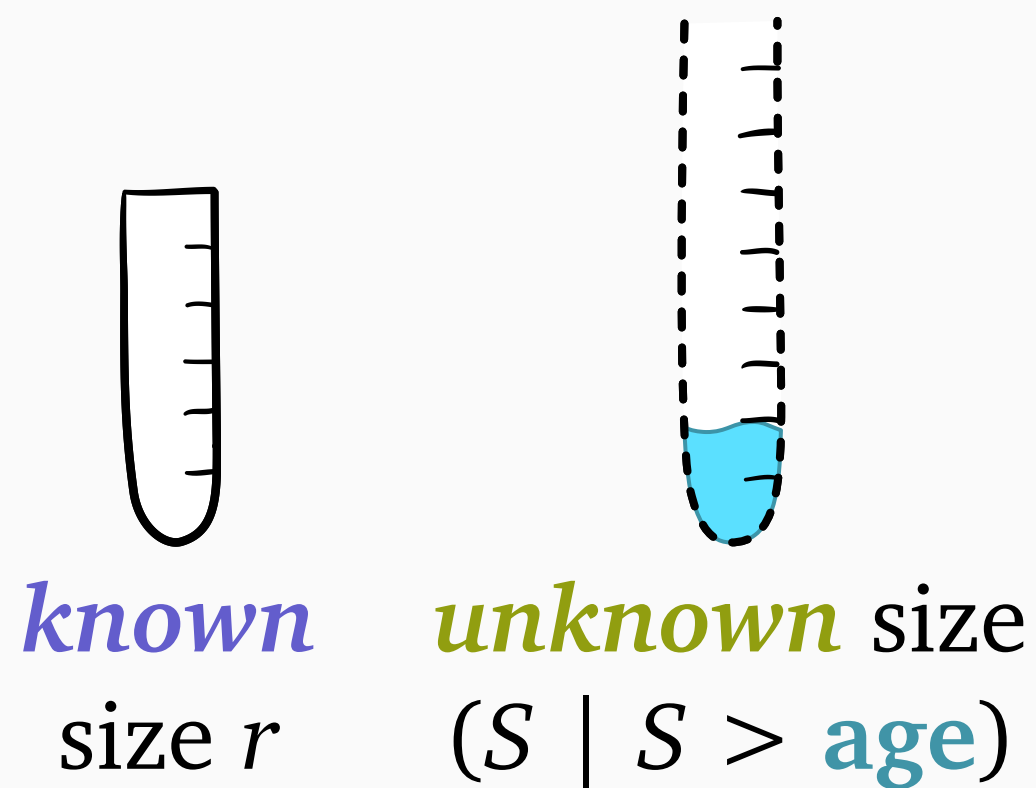
known
size r



unknown size
($S \mid S > \text{age}$)

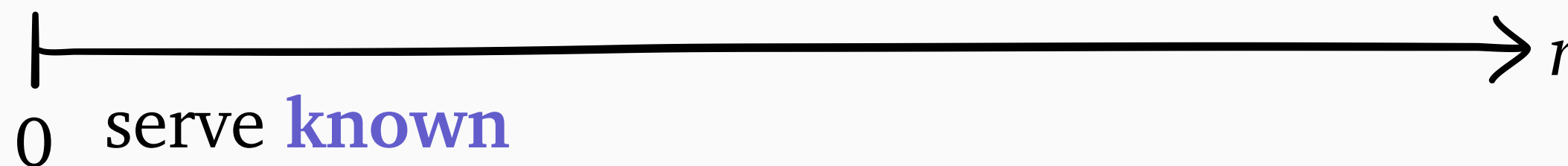
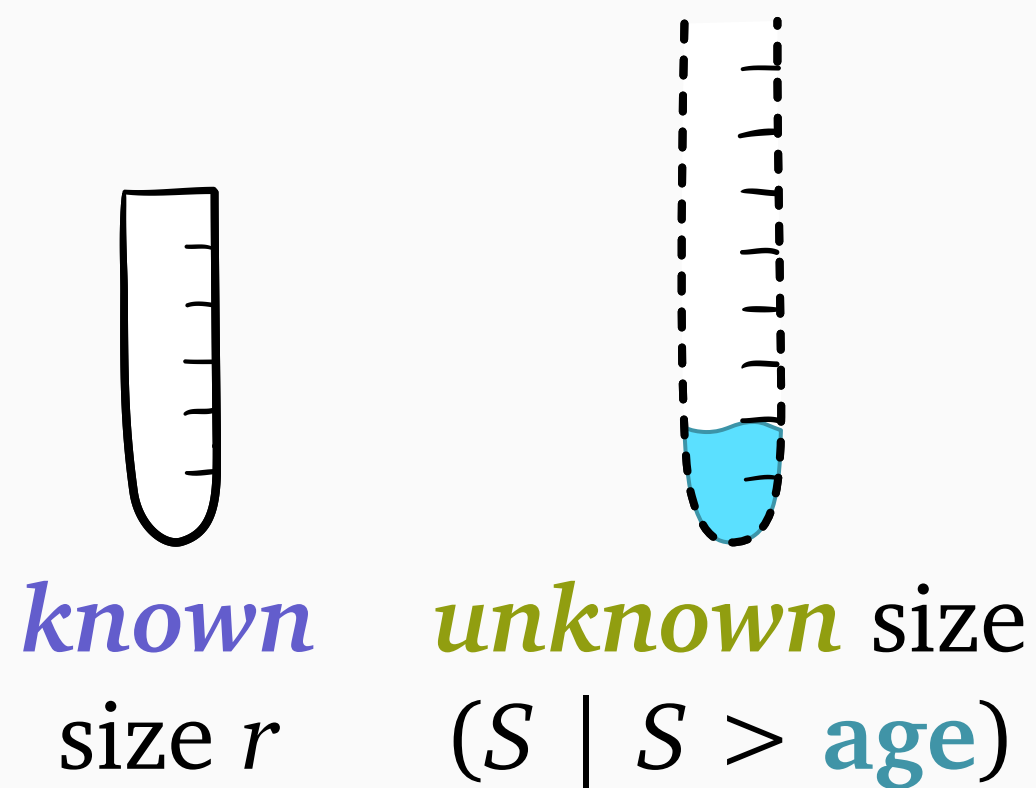


Defining the Gittins rank



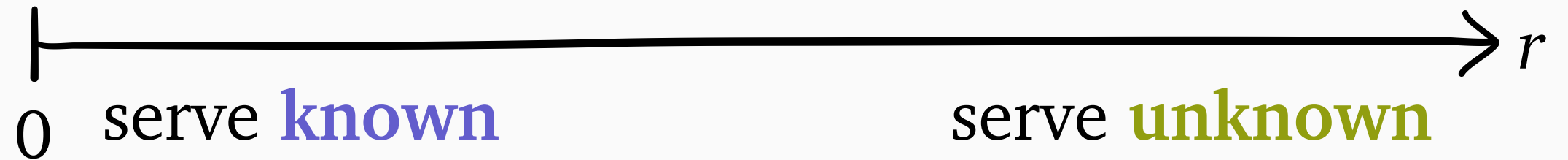
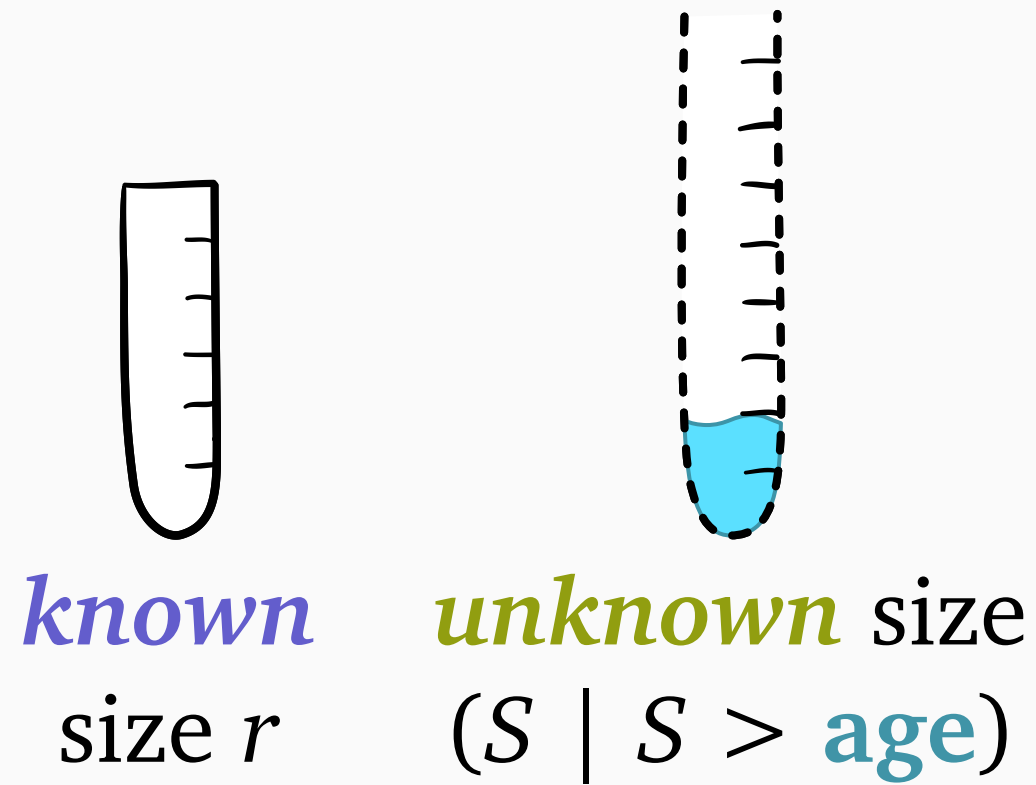
? Key question: for which r should we serve **unknown**?

Defining the Gittins rank



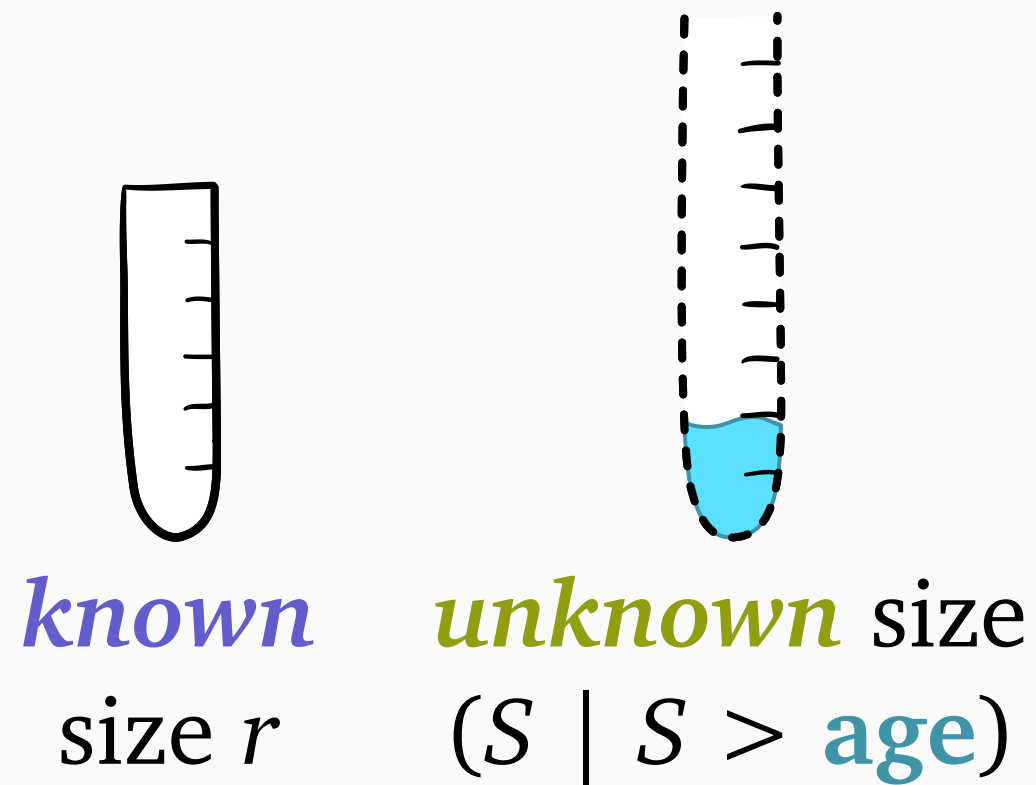
? Key question: for which r should we serve *unknown*?

Defining the Gittins rank



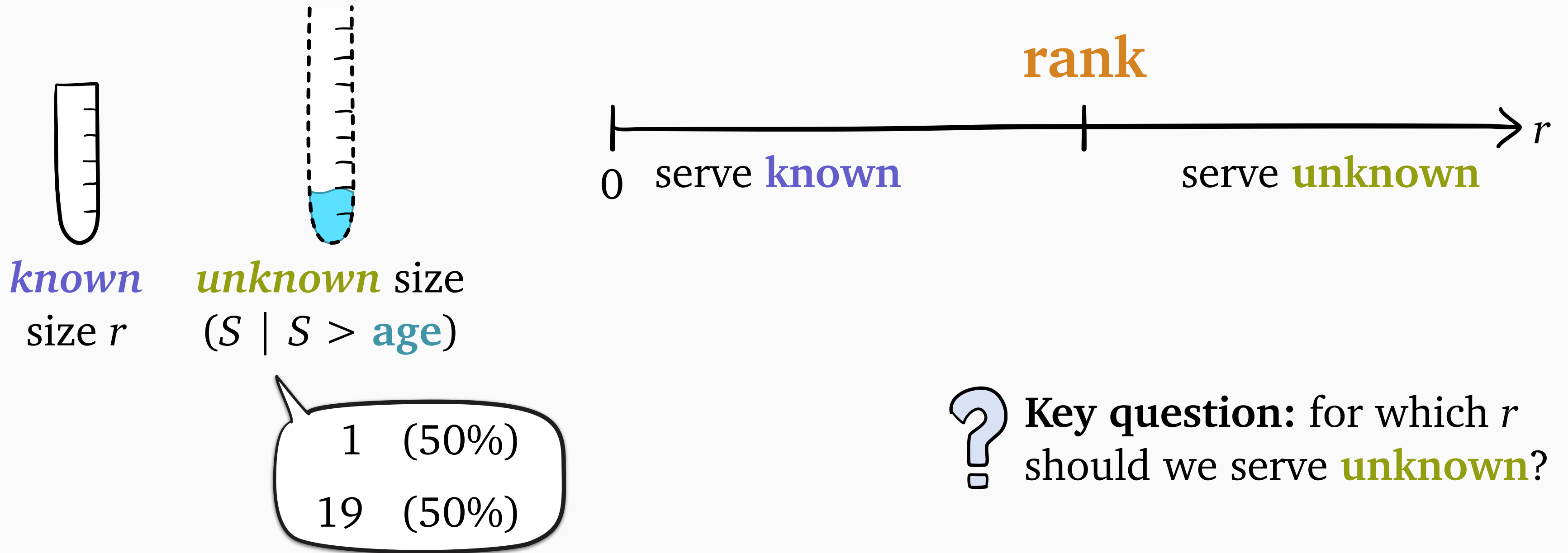
? Key question: for which r should we serve *unknown*?

Defining the Gittins rank



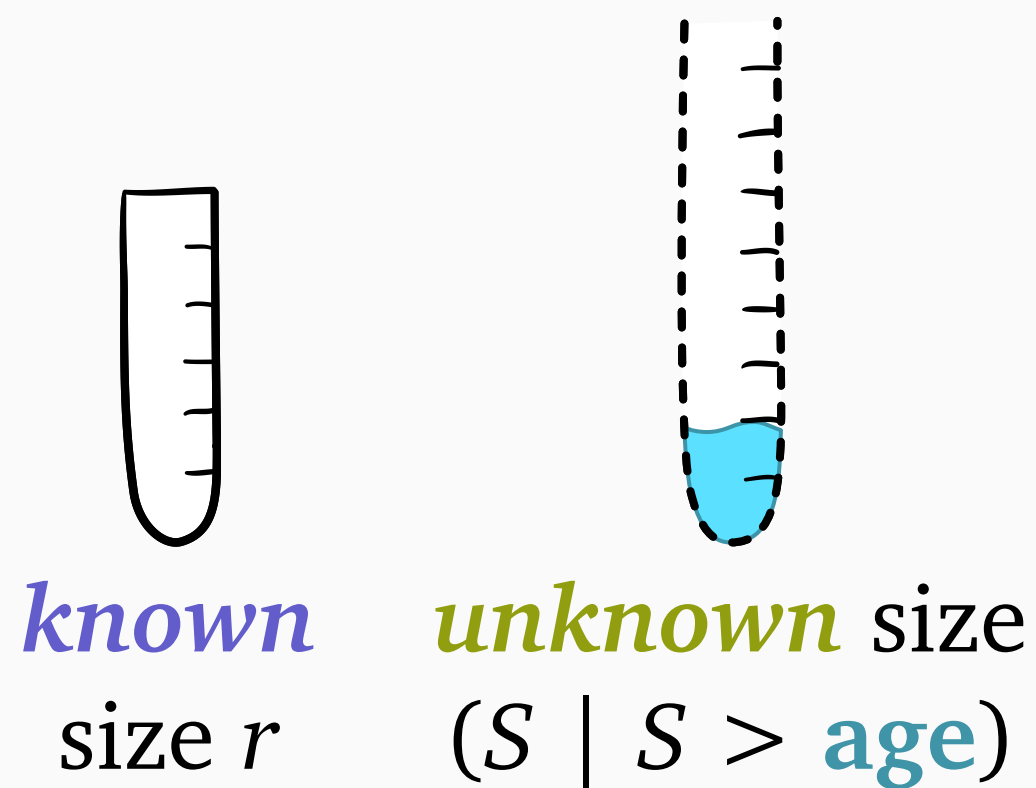
? Key question: for which r should we serve *unknown*?

Defining the Gittins rank



? Key question: for which r should we serve *unknown*?

Defining the Gittins rank

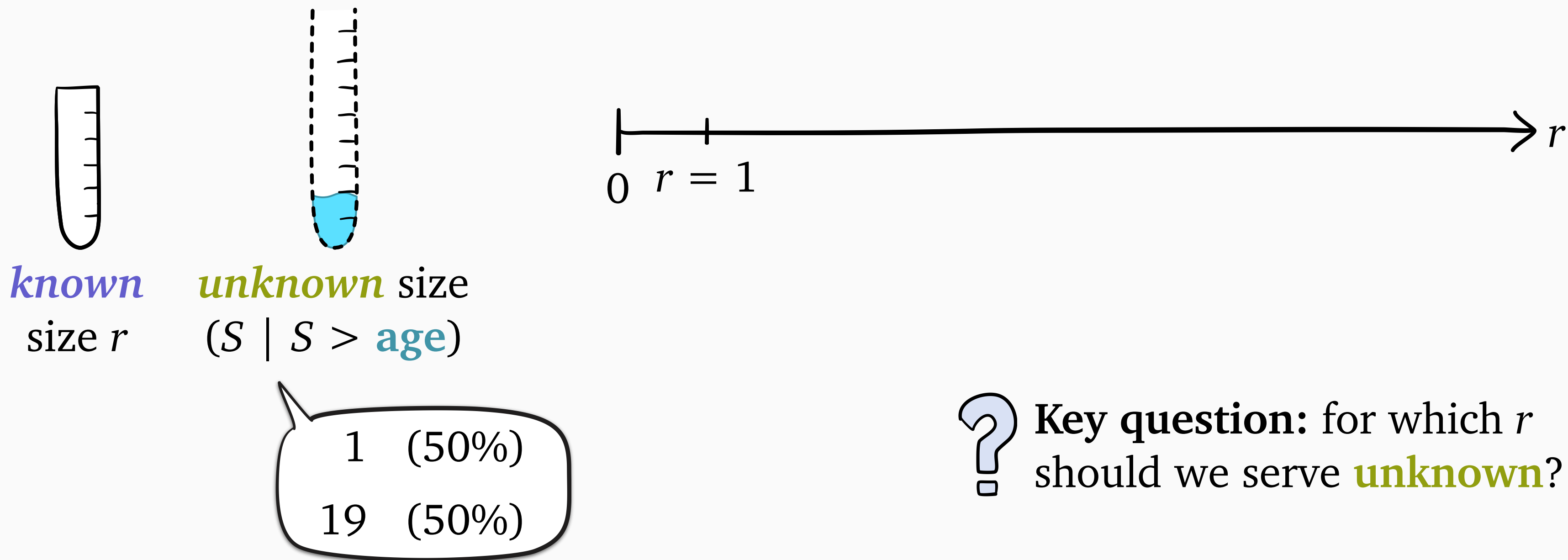


1	(50%)
19	(50%)

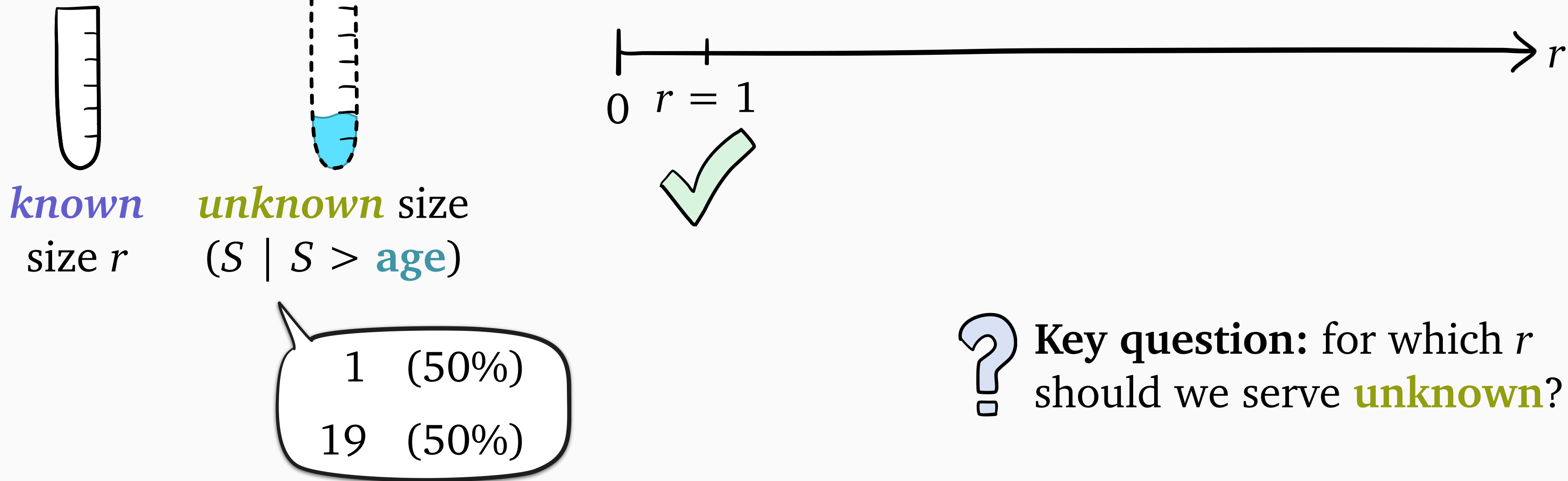


? Key question: for which r should we serve **unknown**?

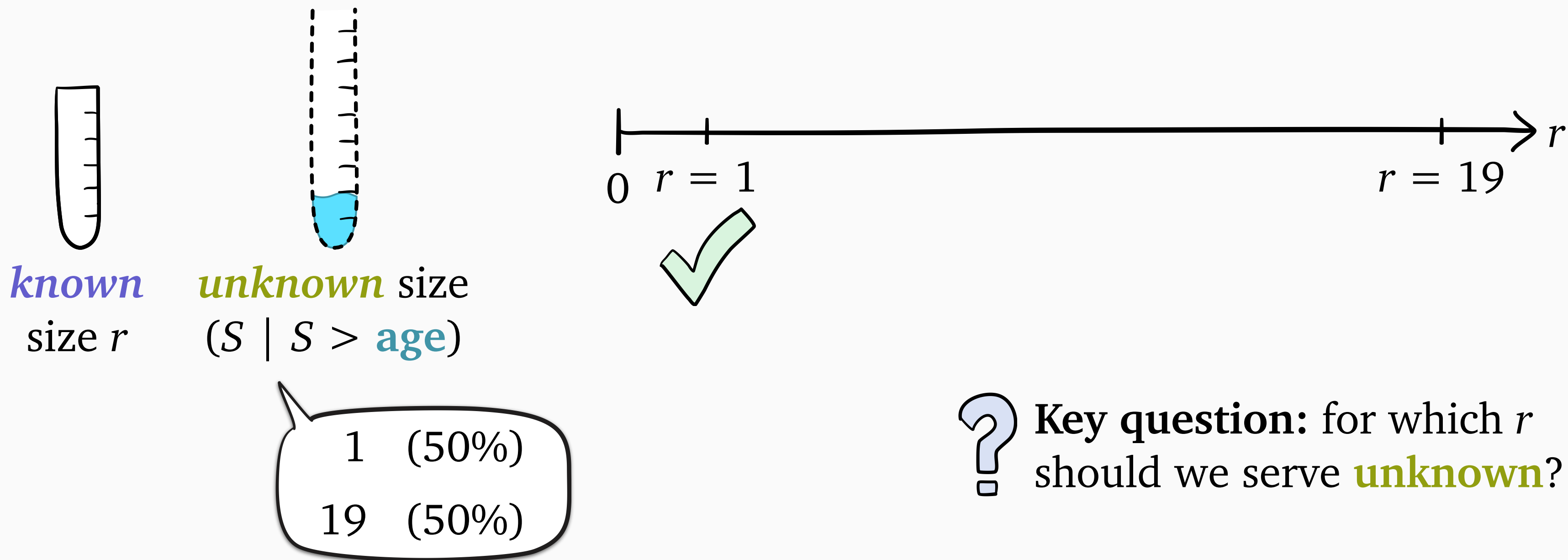
Defining the Gittins rank



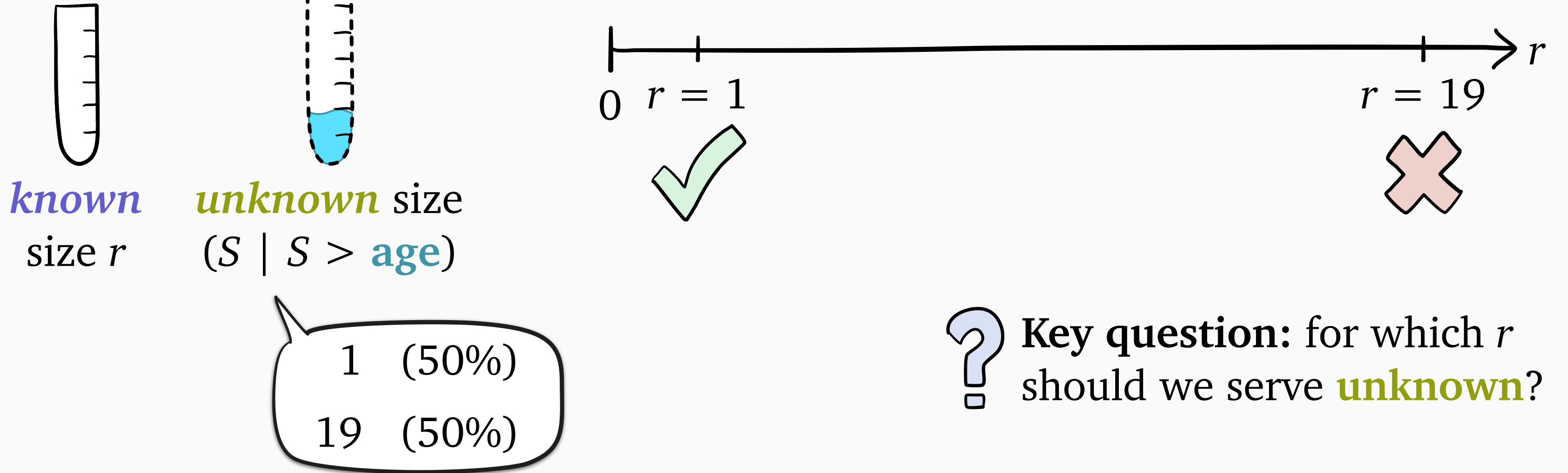
Defining the Gittins rank



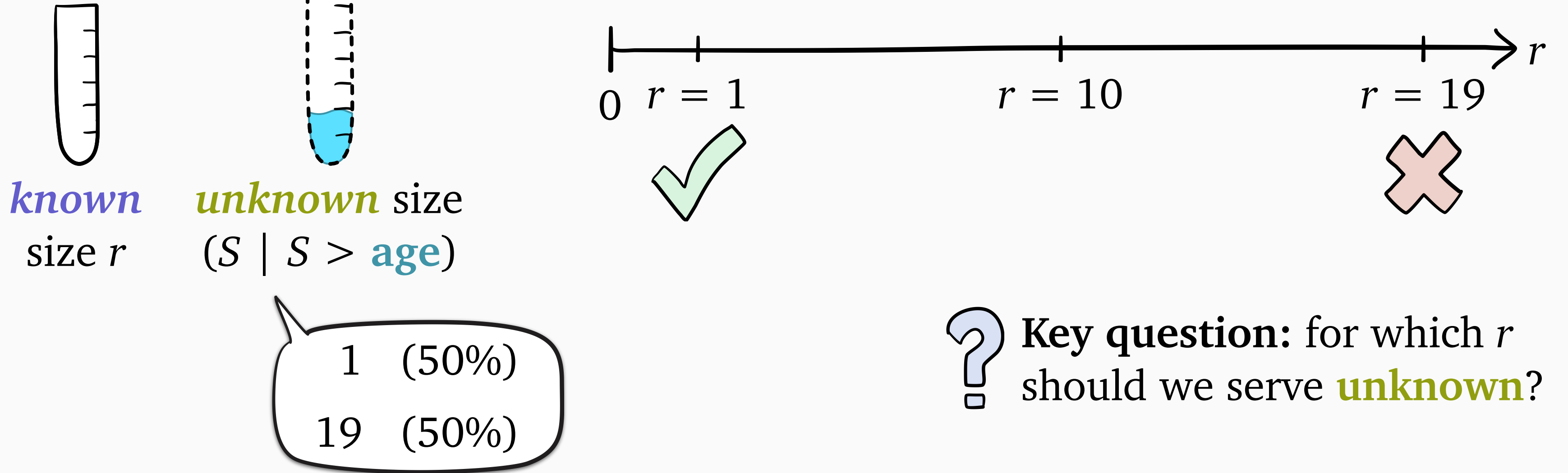
Defining the Gittins rank



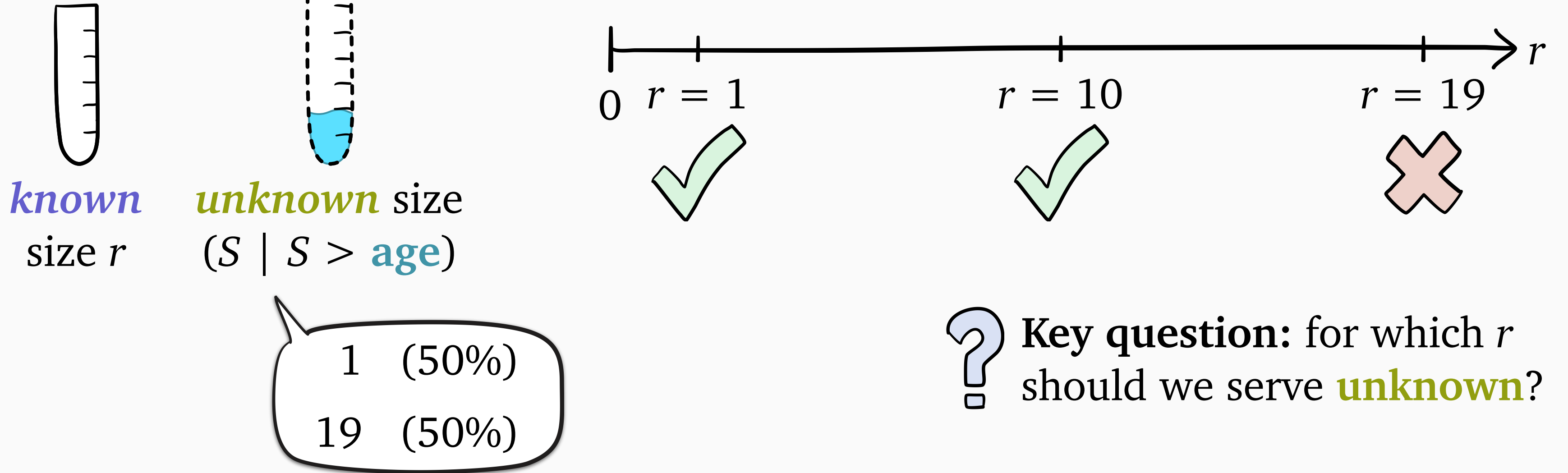
Defining the **Gittins rank**



Defining the Gittins rank

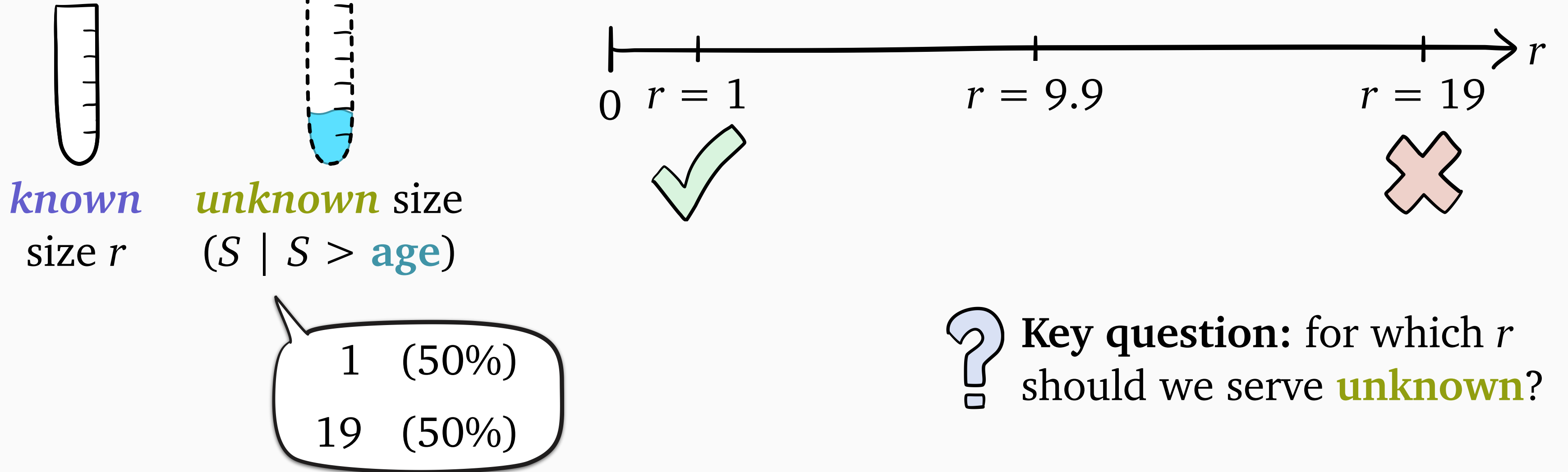


Defining the Gittins rank

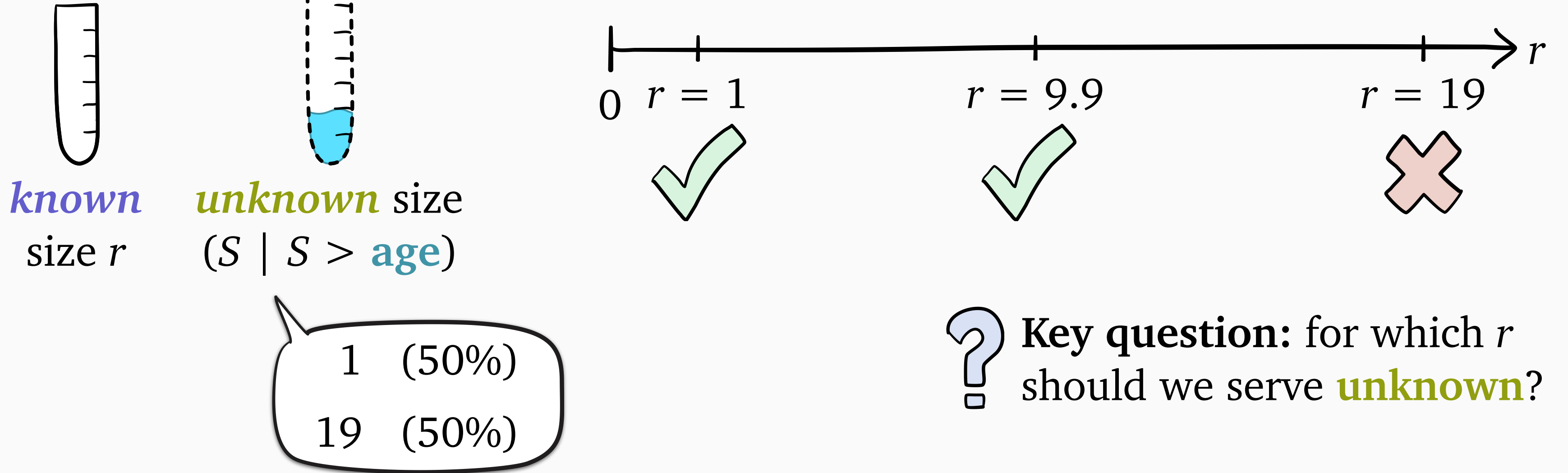


? Key question: for which r should we serve **unknown**?

Defining the Gittins rank

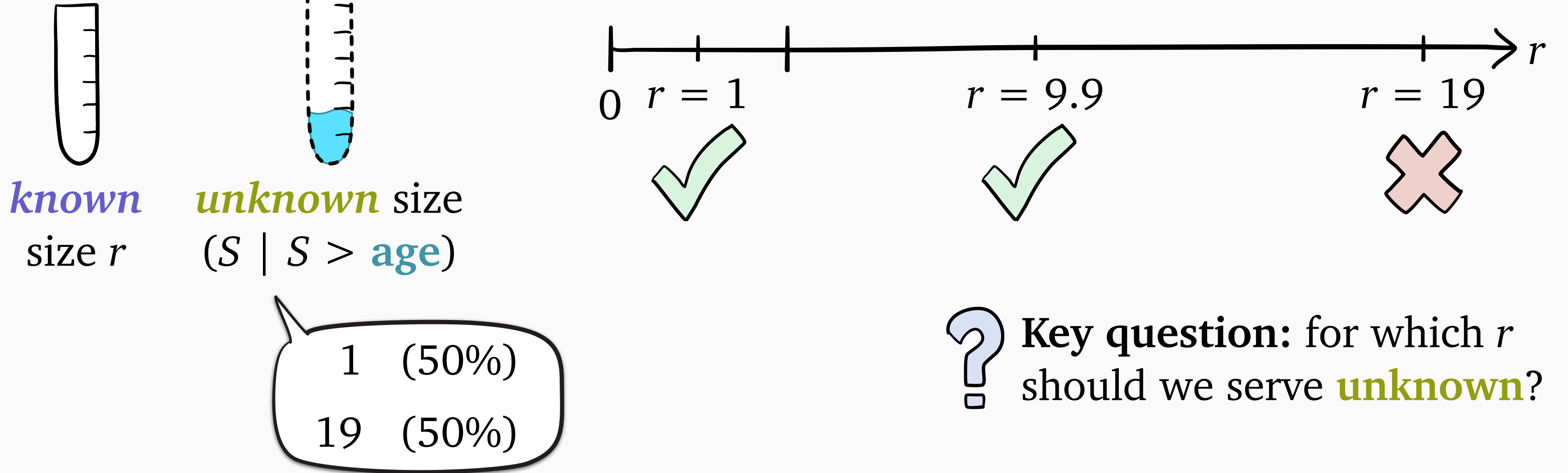


Defining the Gittins rank



? Key question: for which r should we serve **unknown**?

Defining the Gittins rank



? Key question: for which r should we serve **unknown**?

How to prove **SRPT** is optimal?

How to prove **SRPT** is optimal?

Little's law: $E[N] = \lambda E[T]$

How to prove **SRPT** is optimal?

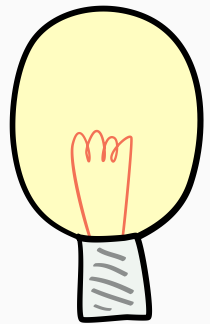
Little's law: $E[N] = \lambda E[T]$

jobs present

How to prove **SRPT** is optimal?

Little's law: $E[N] = \lambda E[T]$

jobs present

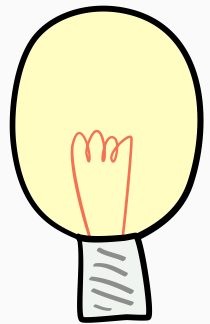


Reward each time
we complete a job

How to prove **SRPT** is optimal?

Little's law: $E[N] = \lambda E[T]$

jobs present



Reward each time
we complete a job

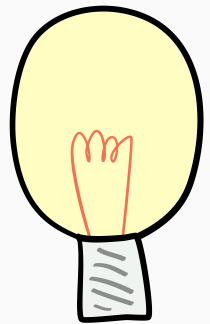


Rewards aren't
immediate

How to prove **SRPT** is optimal?

Little's law: $E[N] = \lambda E[T]$

jobs present



Reward each time
we complete a job



Rewards aren't
immediate

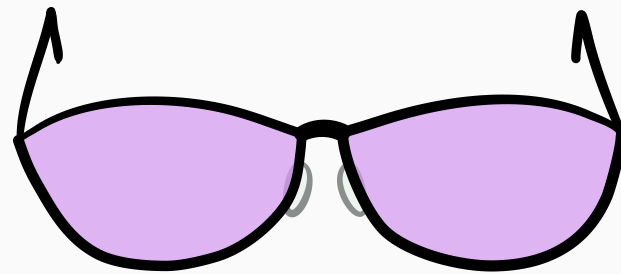


Write N in terms of
“smoother” quantity?

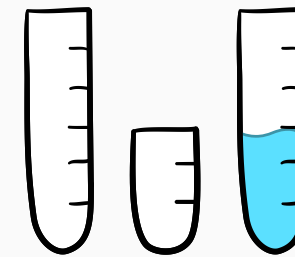
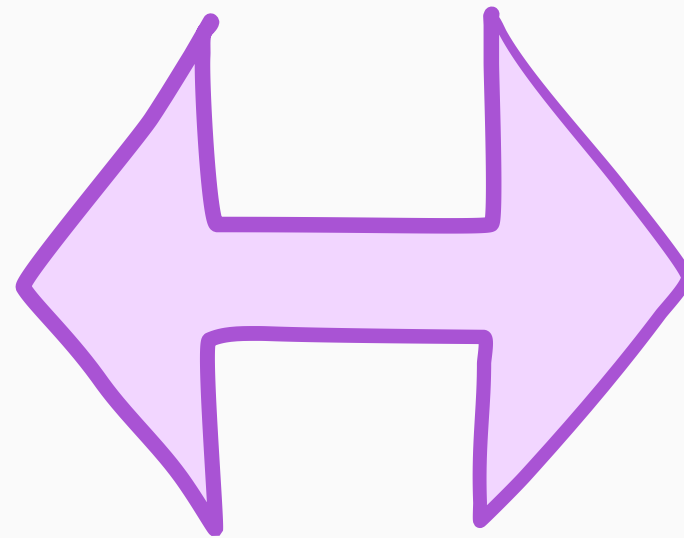


WINE

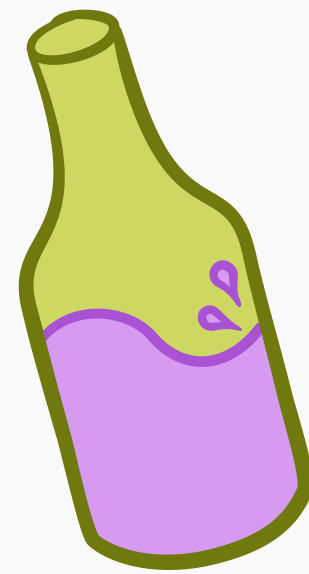
Work Integral Number Equality



r-work $W(r)$

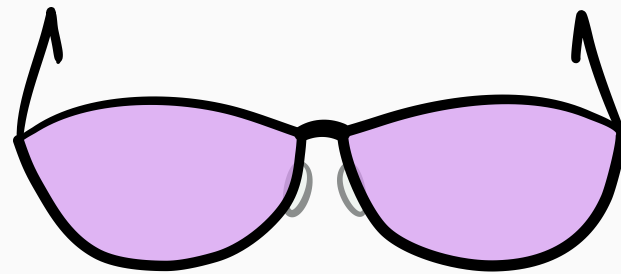


number of jobs N

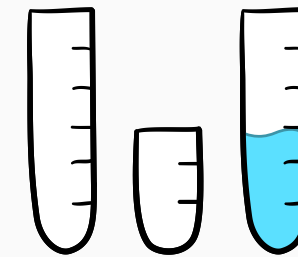
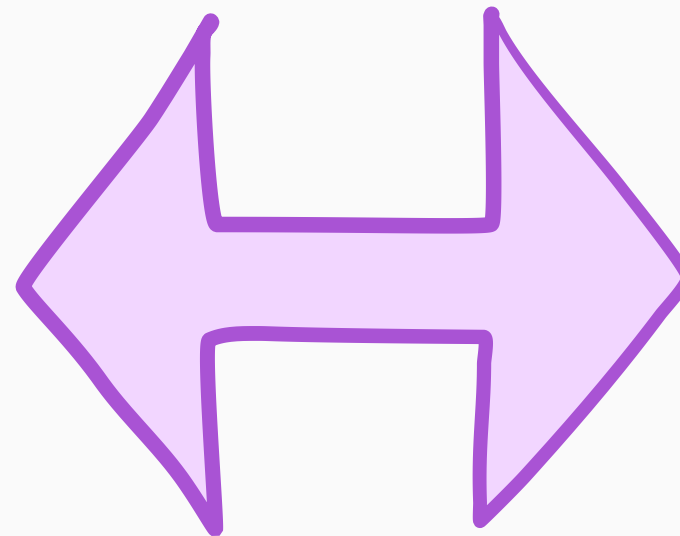


WINE

Work Integral Number Equality



r -work $W(r)$



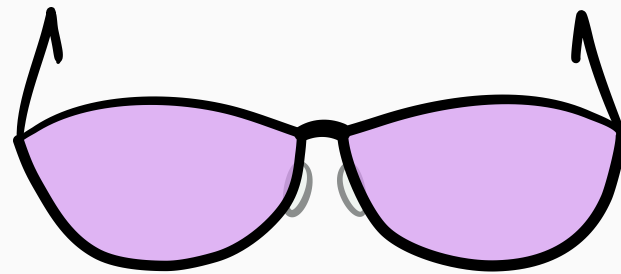
number of jobs N

$E[T]$ bounds for **SRPT** and **Gittins** in:

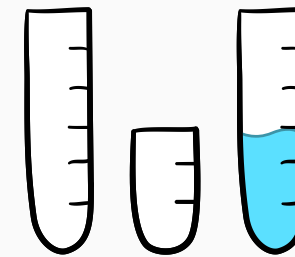
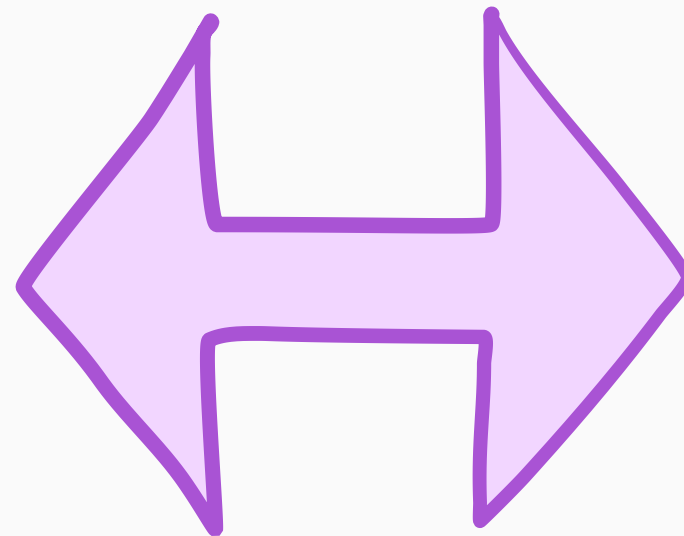


WINE

Work Integral Number Equality



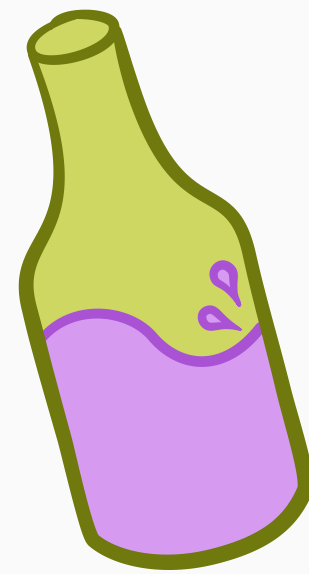
r -work $W(r)$



number of jobs N

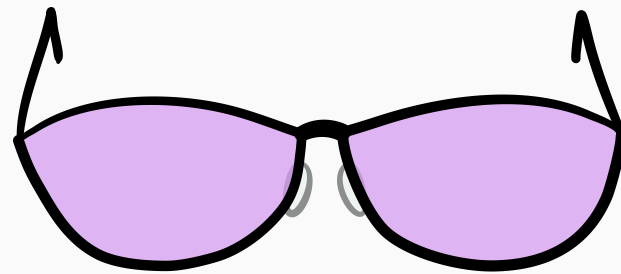
$E[T]$ bounds for **SRPT** and **Gittins** in:

- **M/G/k** and **G/G/k**

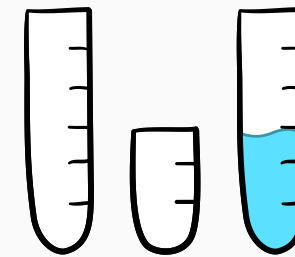
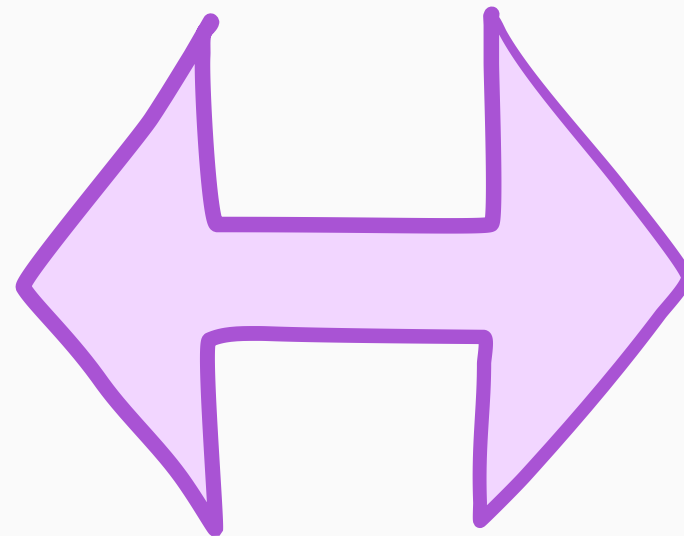


WINE

Work Integral Number Equality



r -work $W(r)$



number of jobs N

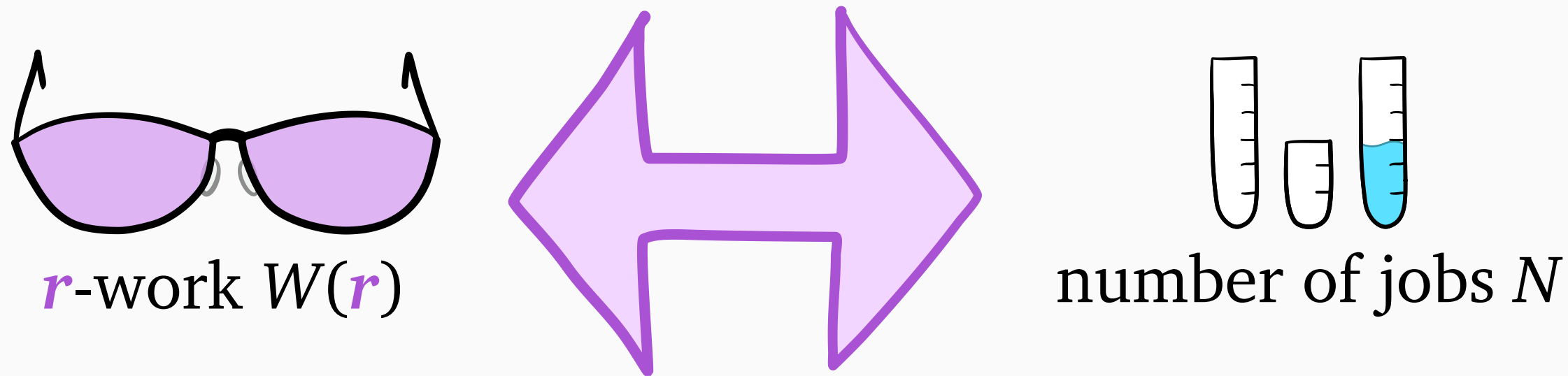
$E[T]$ bounds for **SRPT** and **Gittins** in:

- **M/G/k** and **G/G/k**
- systems with **multiserver jobs**



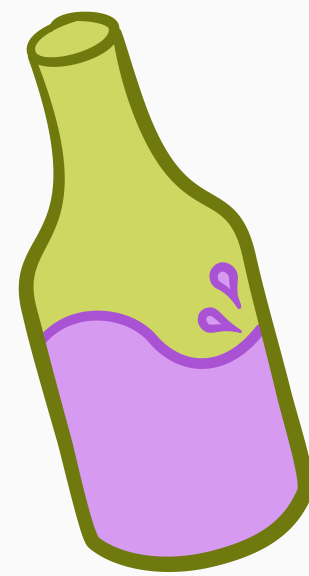
WINE

Work Integral Number Equality



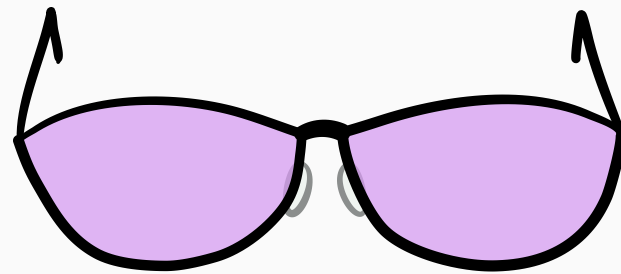
$E[T]$ bounds for **SRPT** and **Gittins** in:

- **M/G/k** and **G/G/k**
- systems with **multiserver jobs**
- systems with **noisy size estimates**

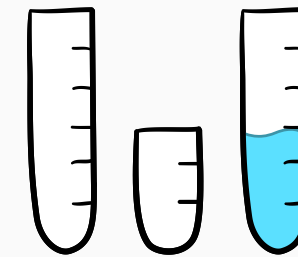
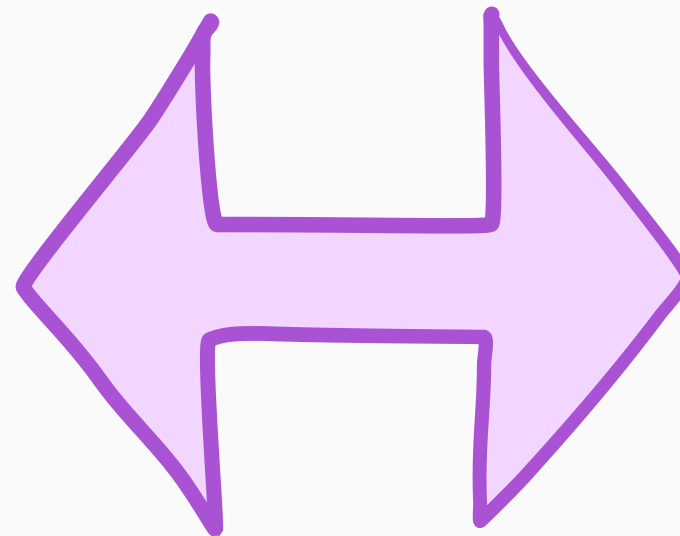


WINE

Work Integral Number Equality



r -work $W(r)$



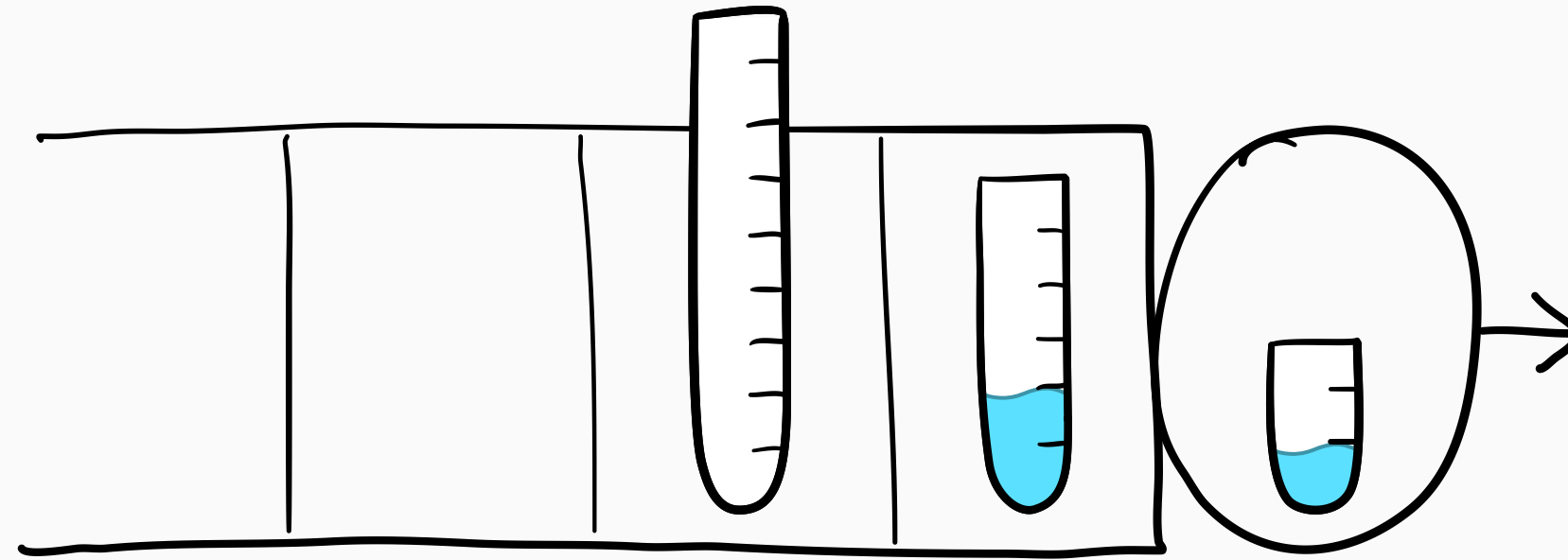
number of jobs N

$E[T]$ bounds for **SRPT** and **Gittins** in:

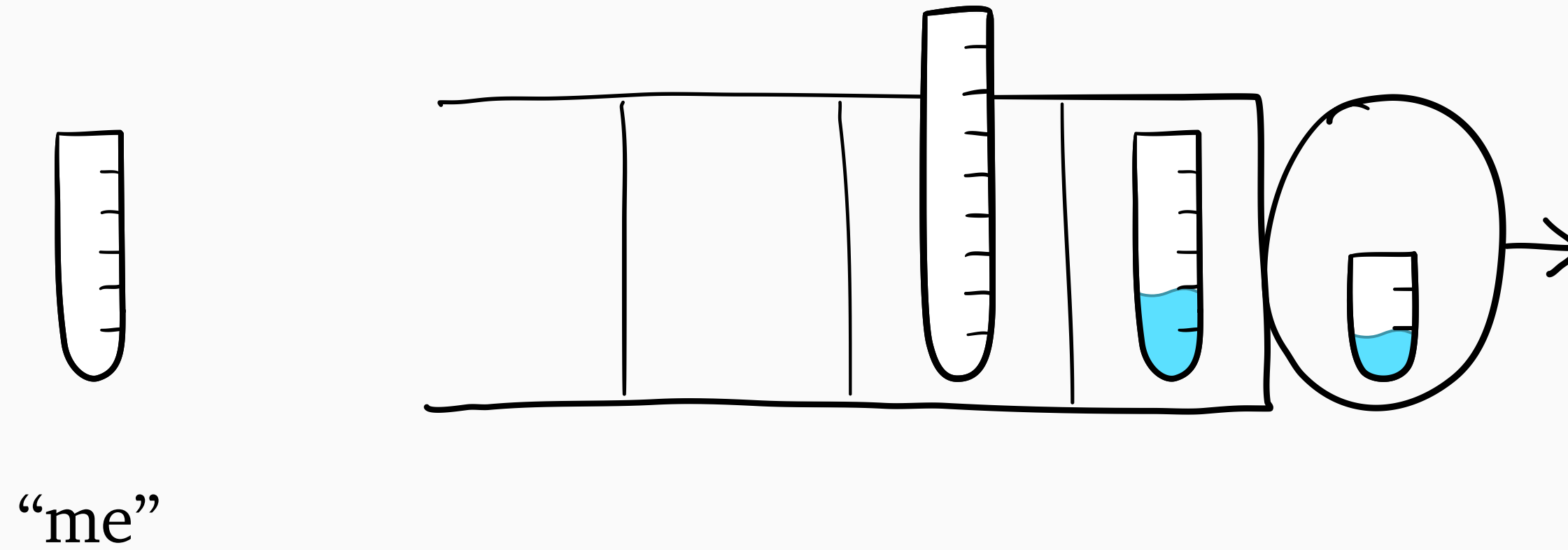
- **M/G/k** and **G/G/k**
- systems with **multiserver jobs**
- systems with **noisy size estimates**
- systems with **unknown size distribution**



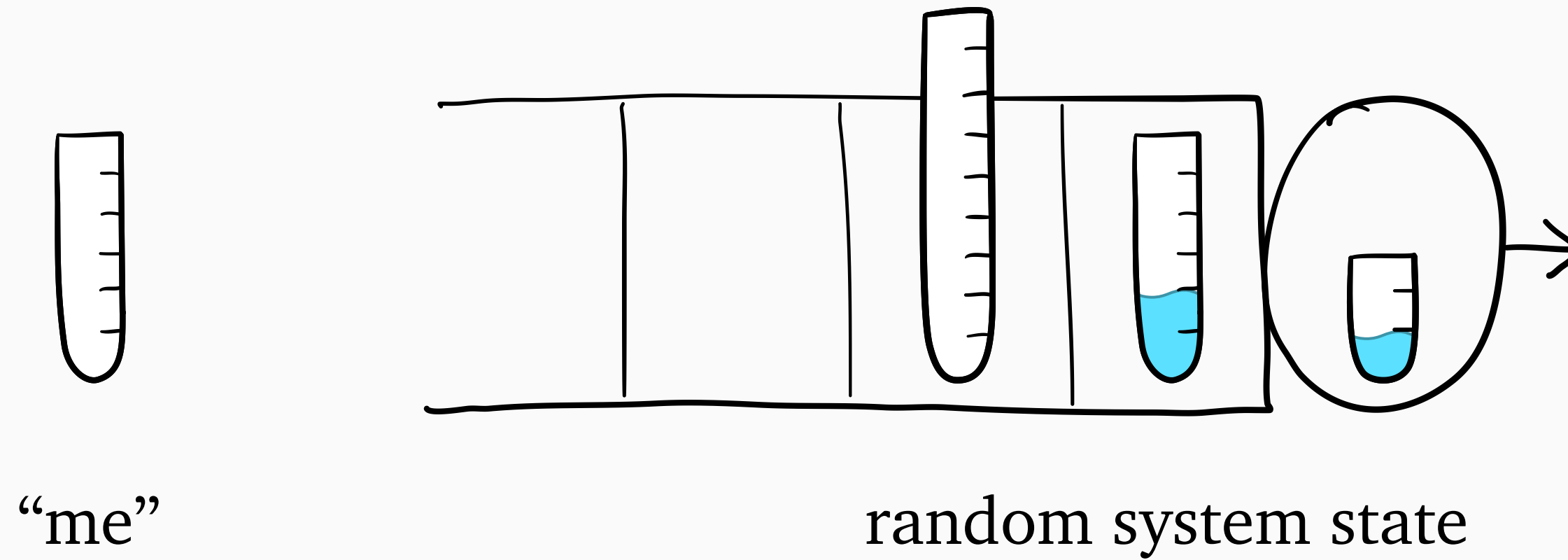
Key quantity: *r*-work



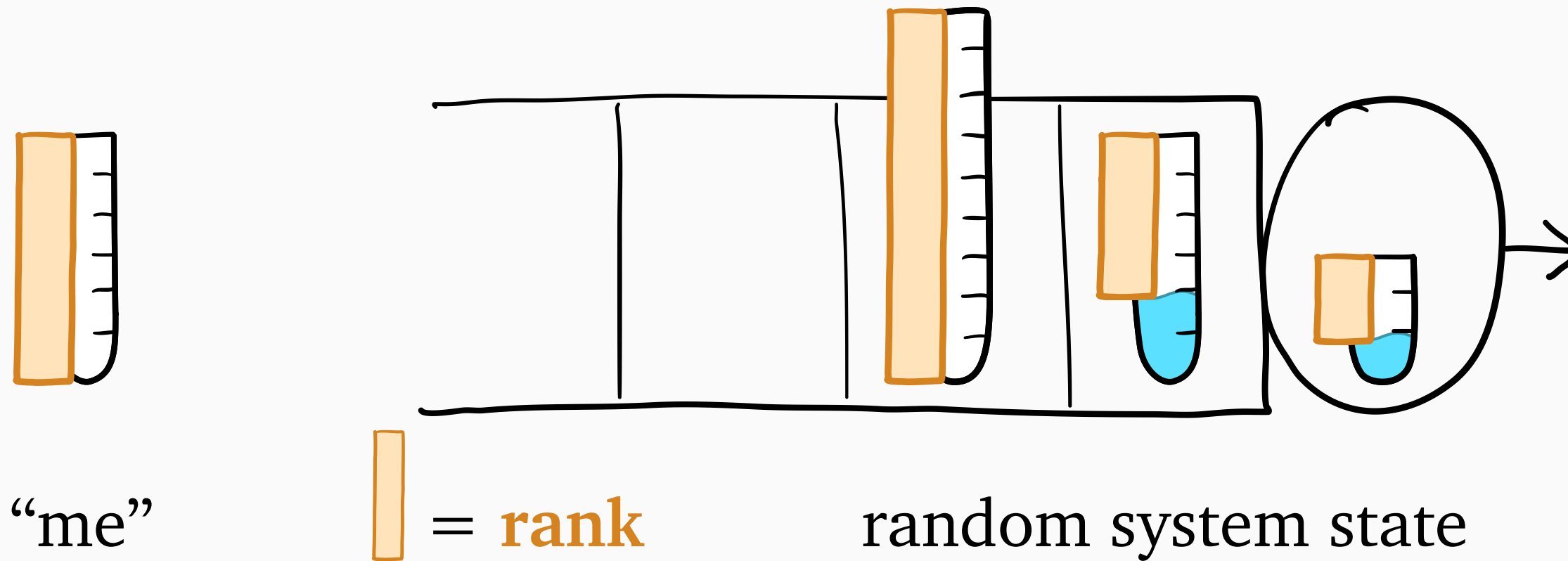
Key quantity: *r*-work



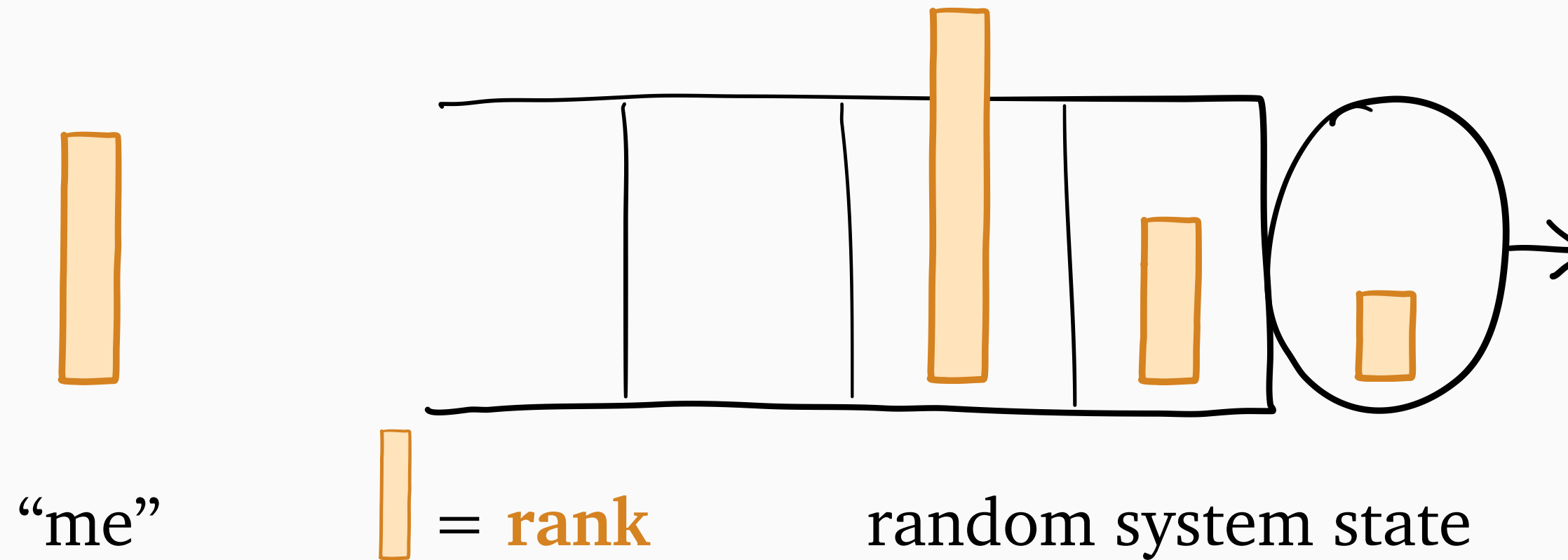
Key quantity: r -work



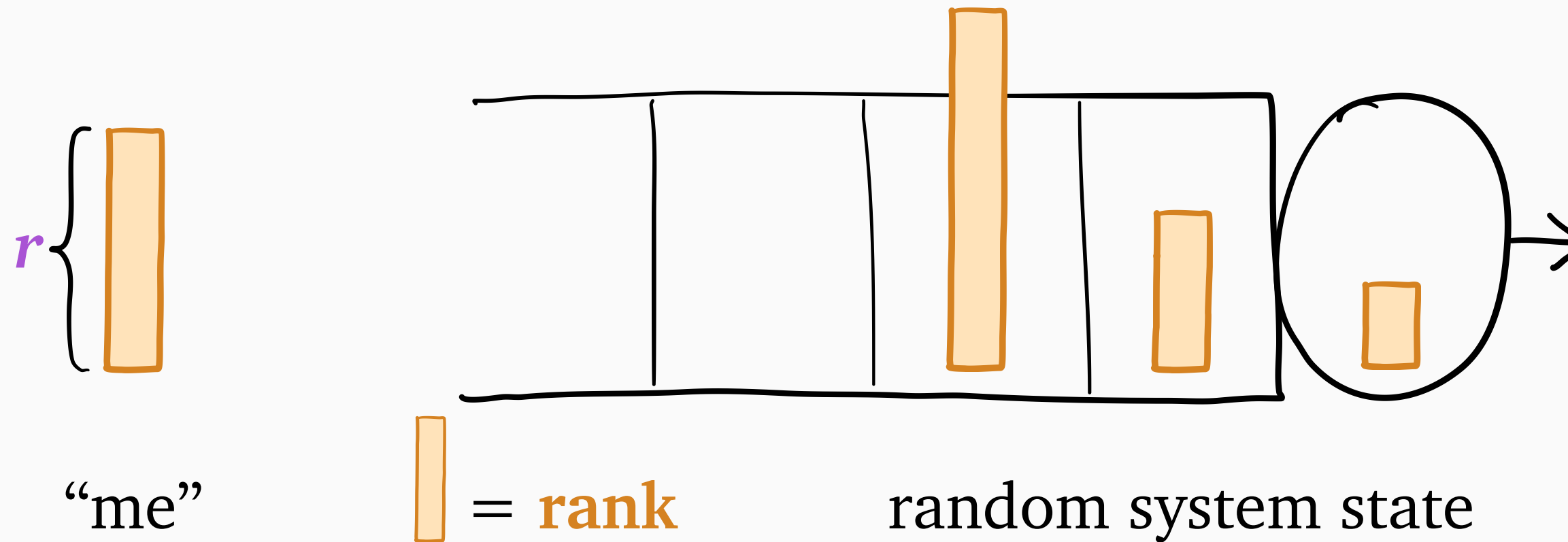
Key quantity: r -work



Key quantity: r -work



Key quantity: r -work

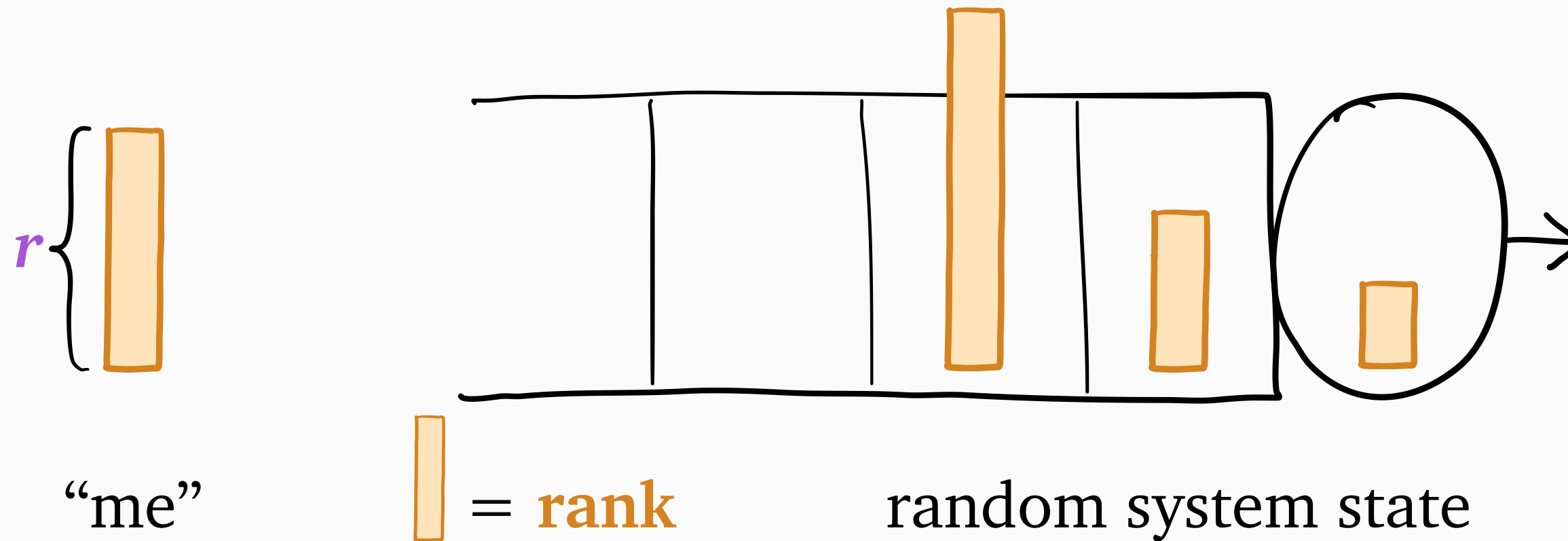


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work

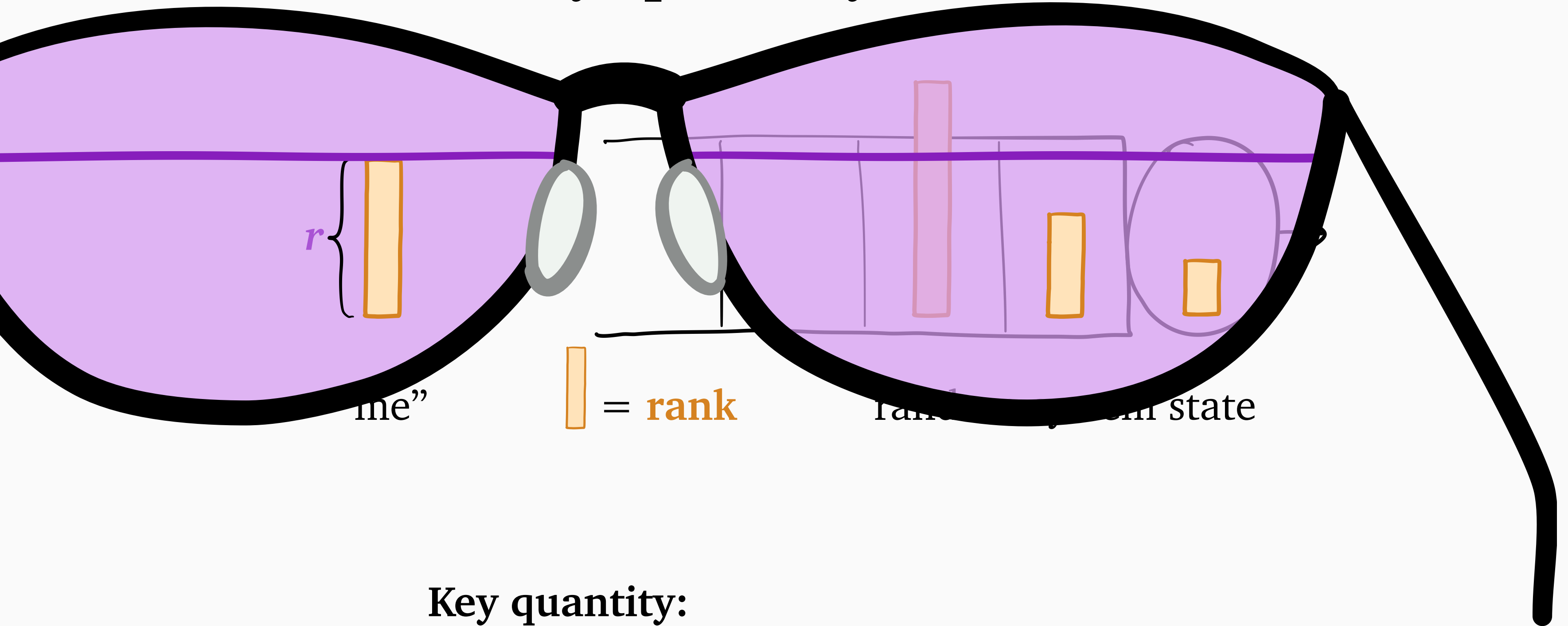


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work

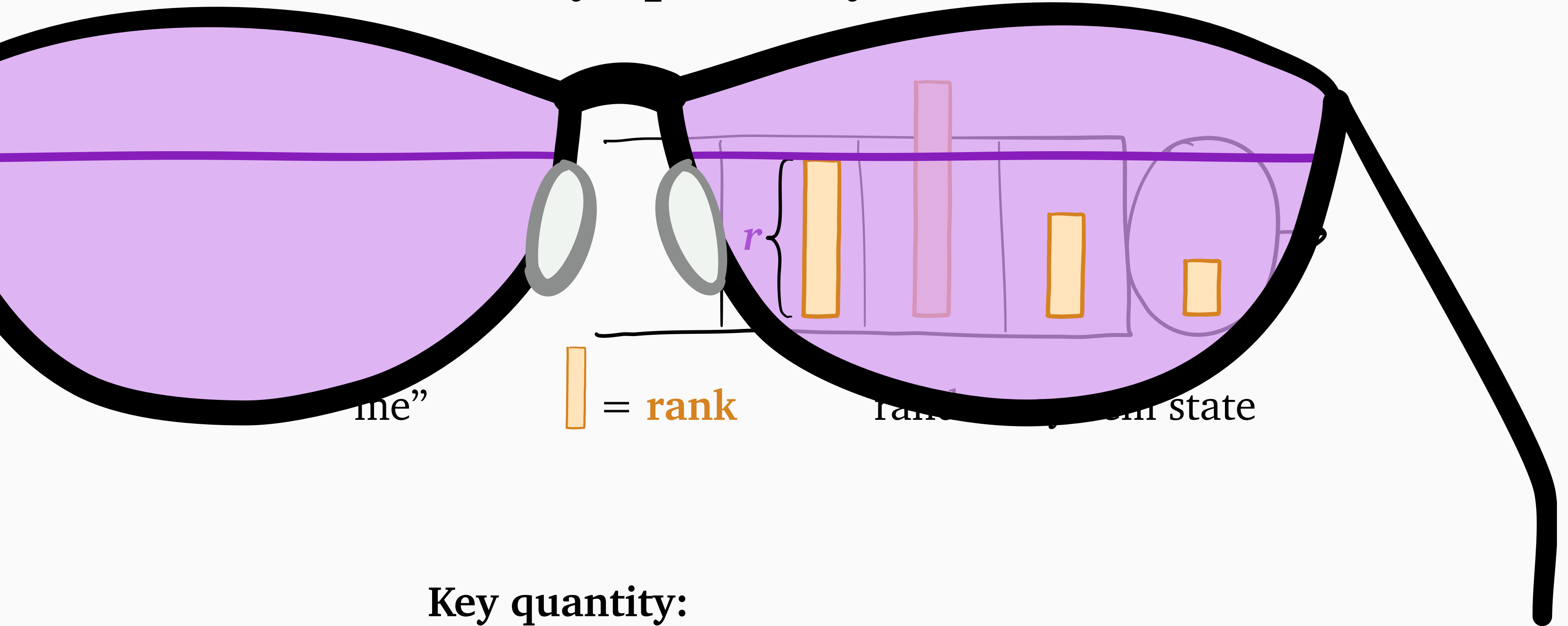


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work

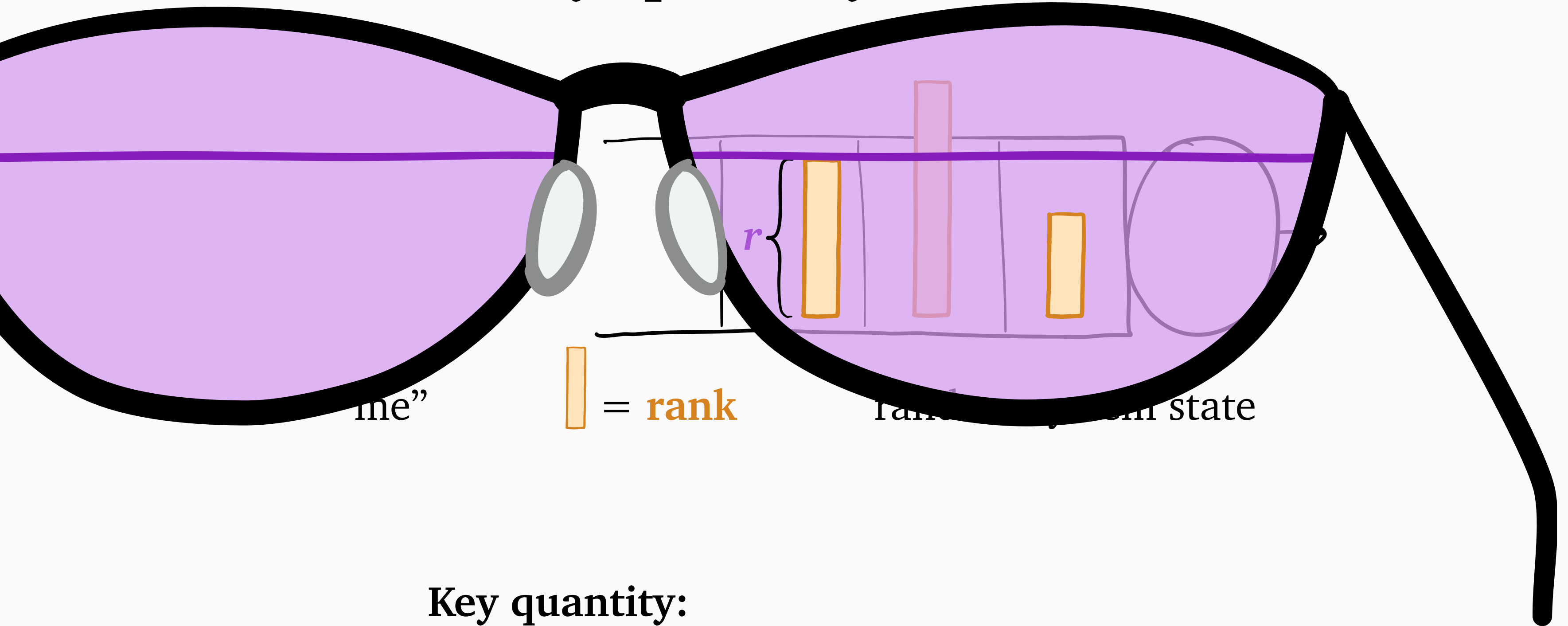


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work

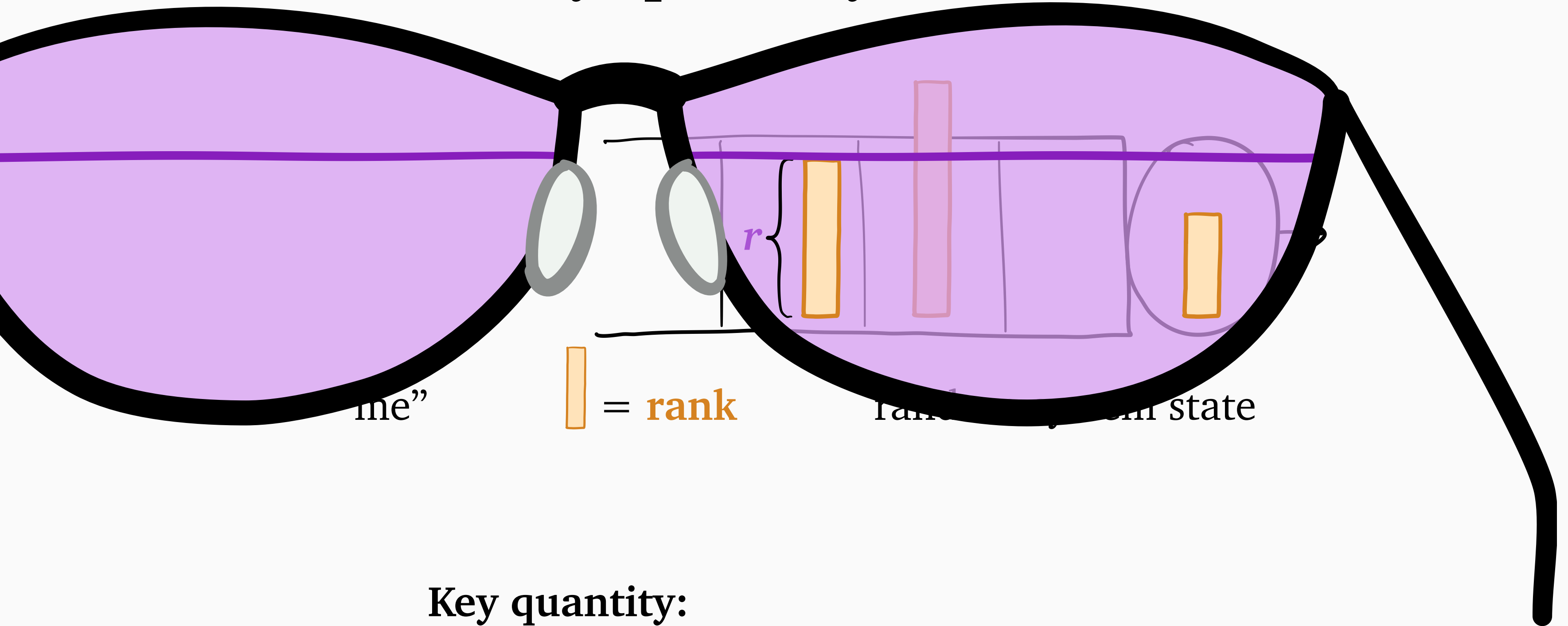


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work

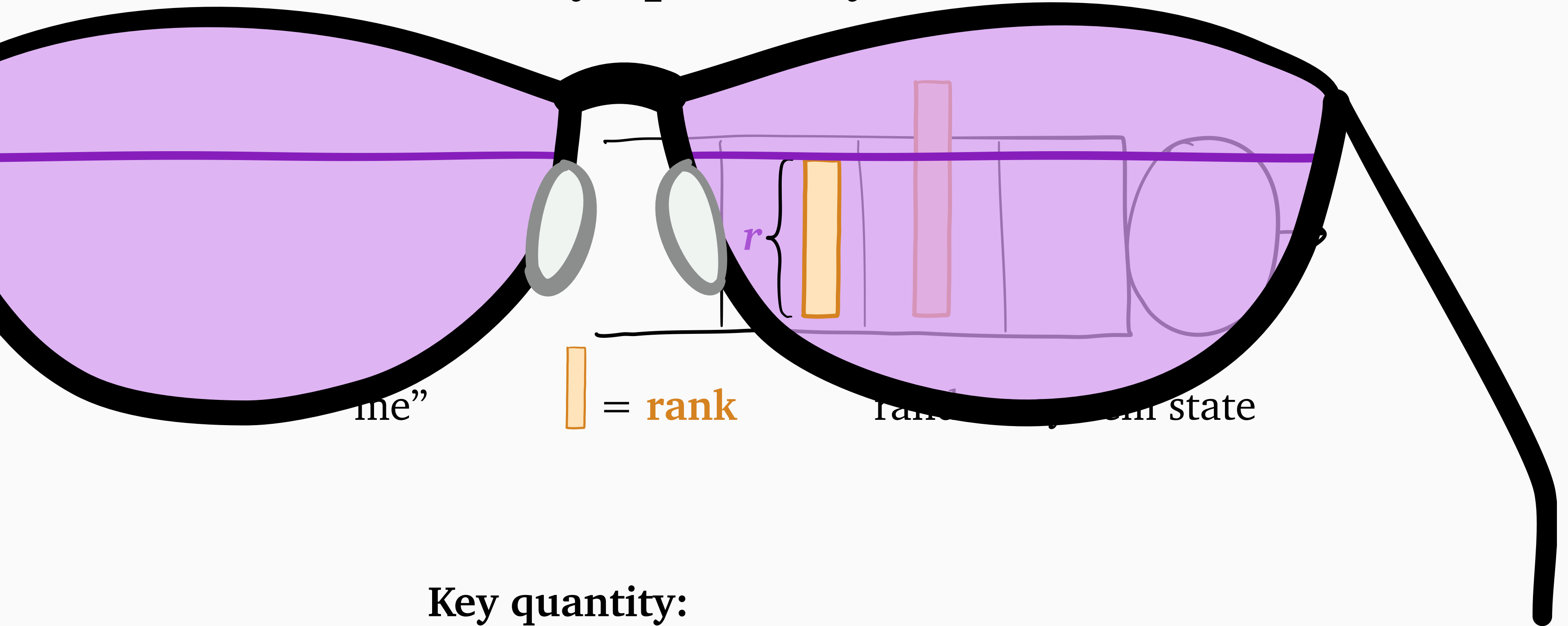


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work

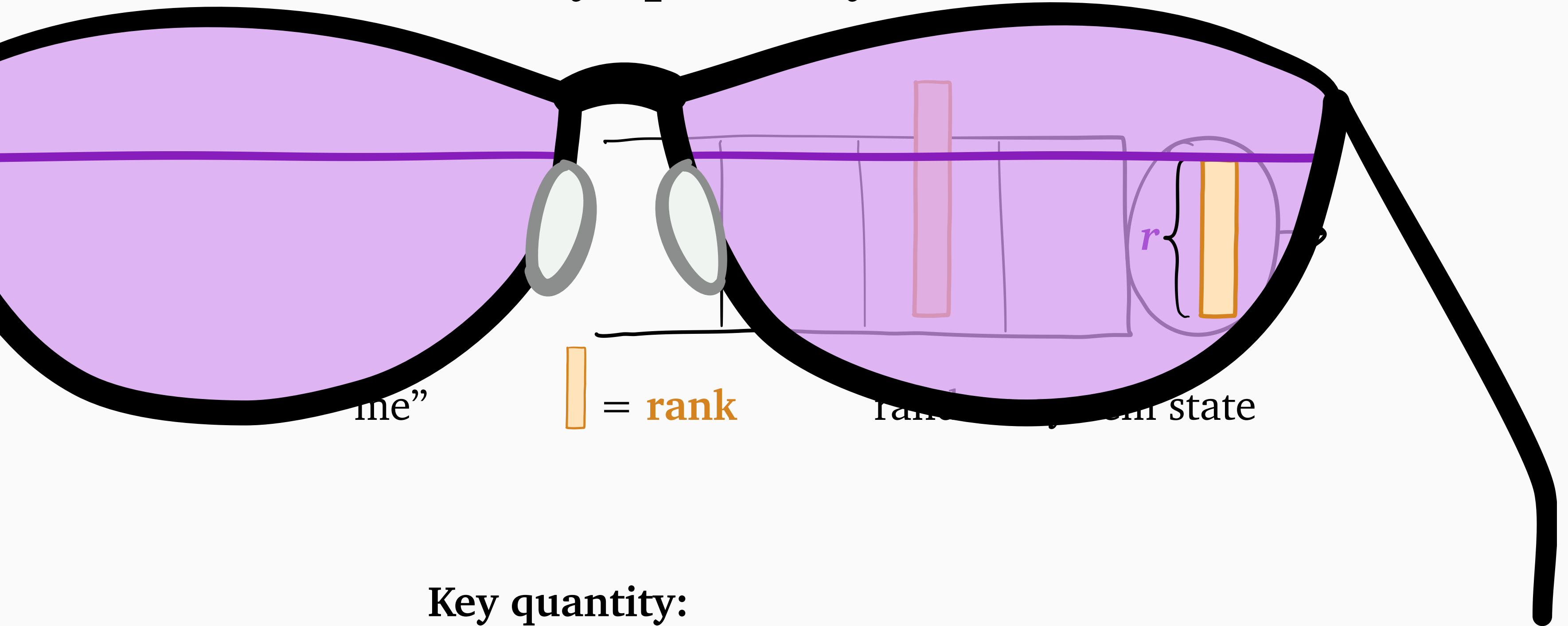


Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Key quantity: r -work



Key quantity:

$W(r)$ = work relevant to job of rank r

r -work

Defining one job's *r*-work

$W(r)$ = work relevant to rank *r*

Defining one job's r -work

$W(r)$ = work relevant to **rank** r

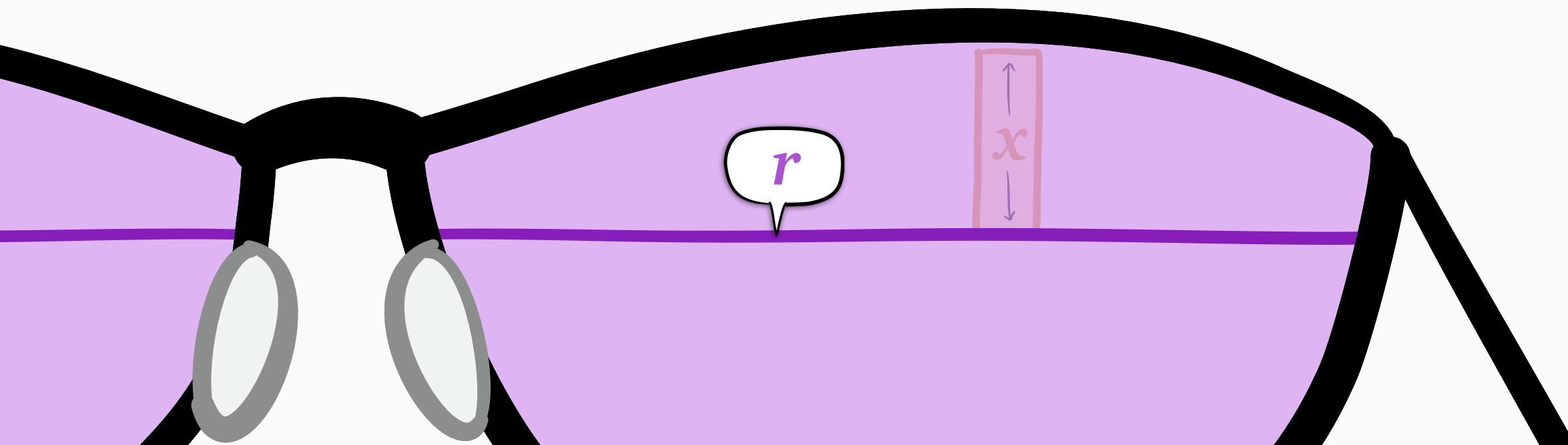
$w_x(r)$ = r -work of *single job* of rem. size x = {



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

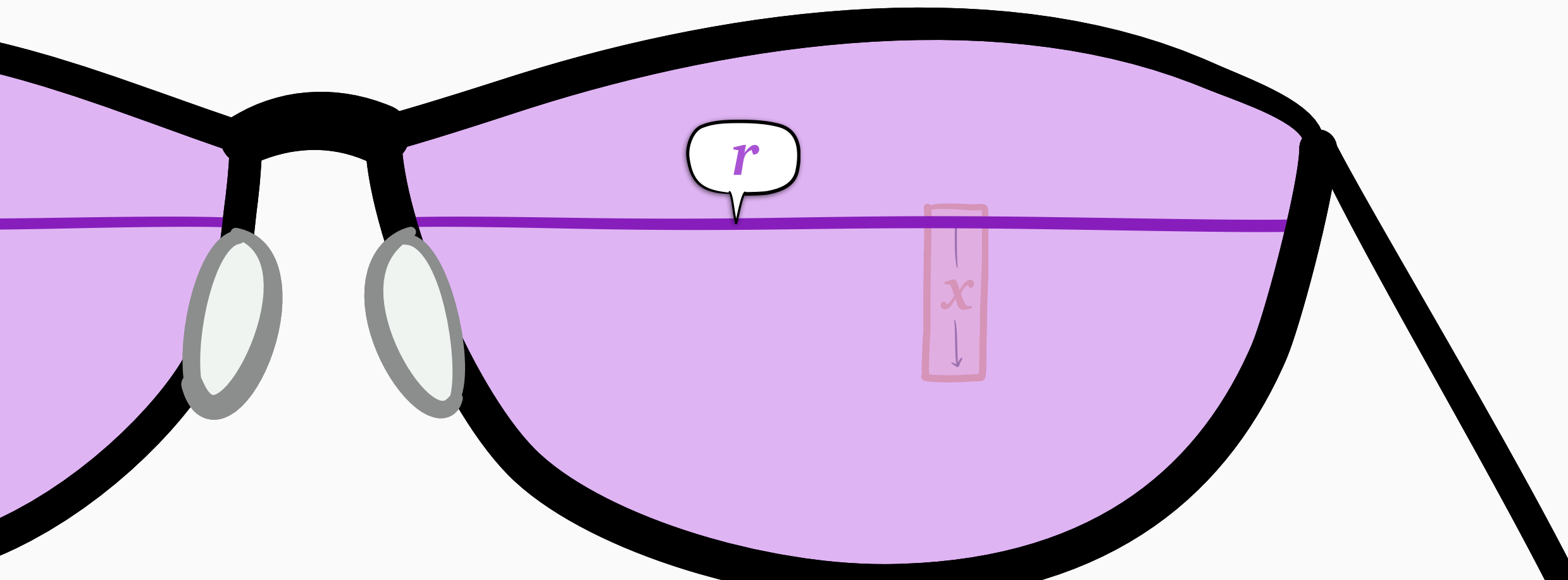
$w_x(r)$ = r -work of *single job* of rem. size x = {



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

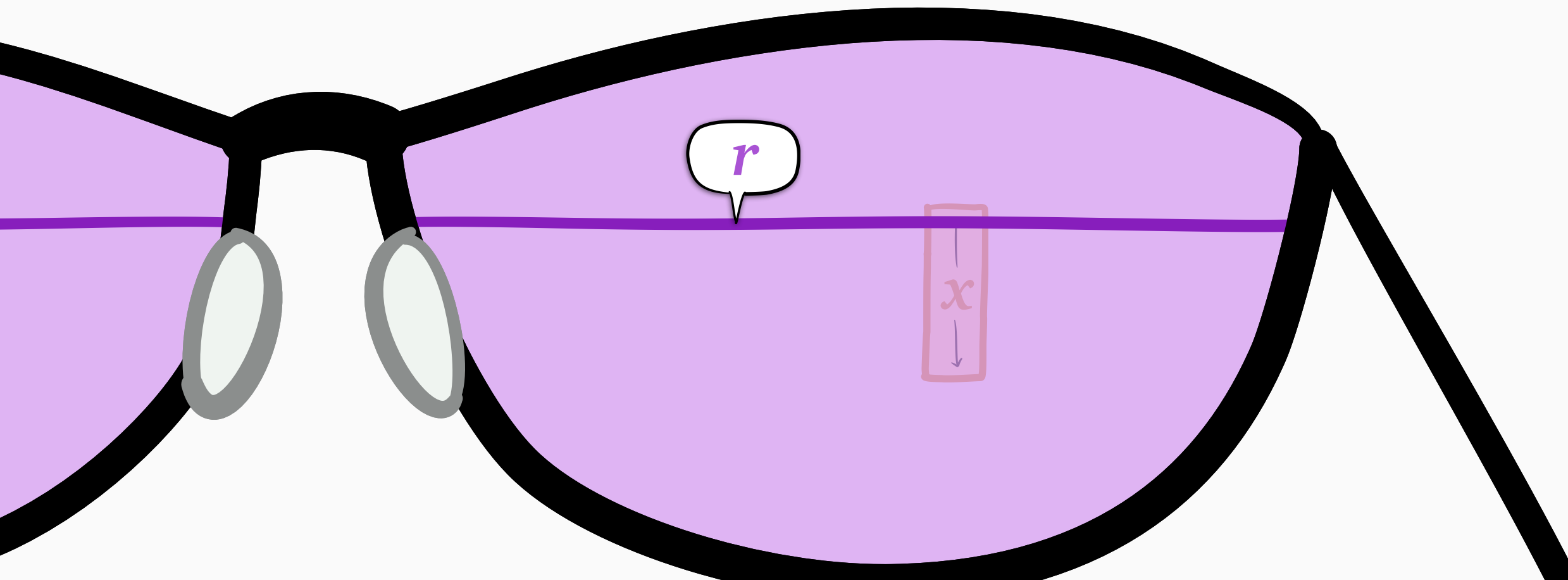
$w_x(r)$ = r -work of *single job* of rem. size x = {



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

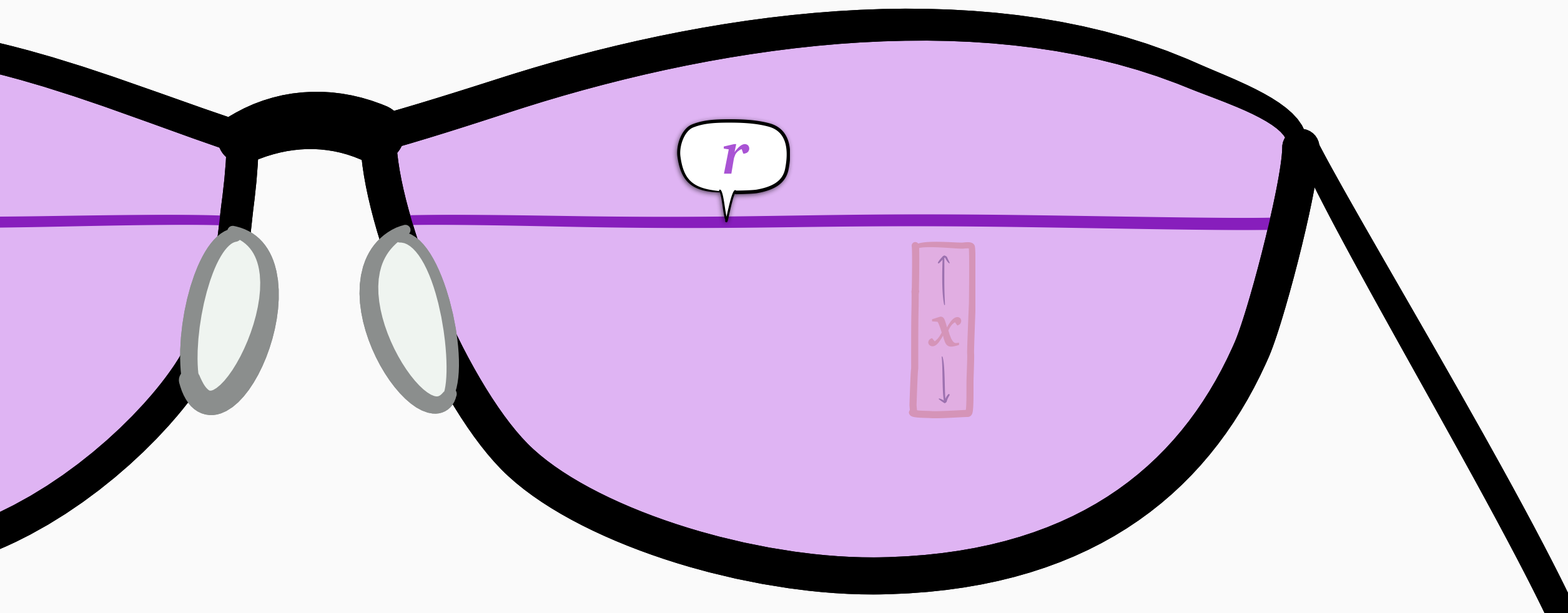
$w_x(r)$ = r -work of *single job* of rem. size x = $\begin{cases} 0 & \text{if } r < x \end{cases}$



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

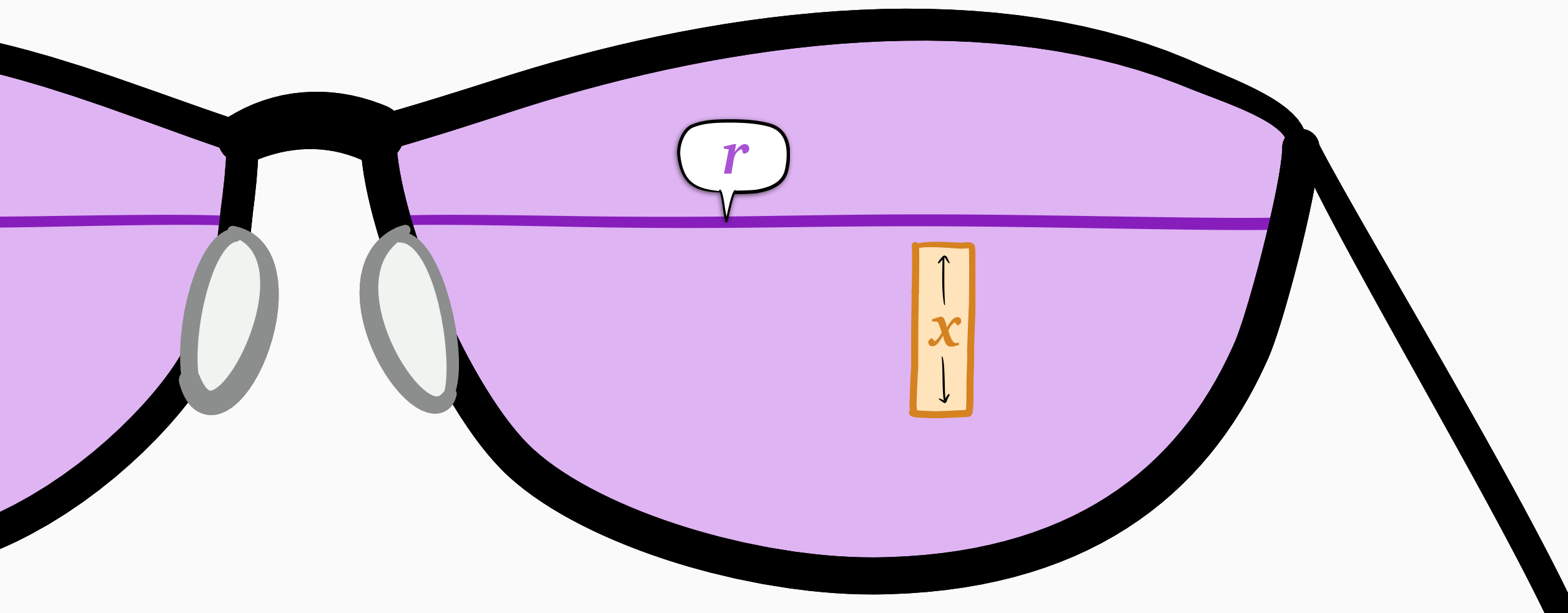
$w_x(r)$ = r -work of *single job* of rem. size x = $\begin{cases} 0 & \text{if } r < x \end{cases}$



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

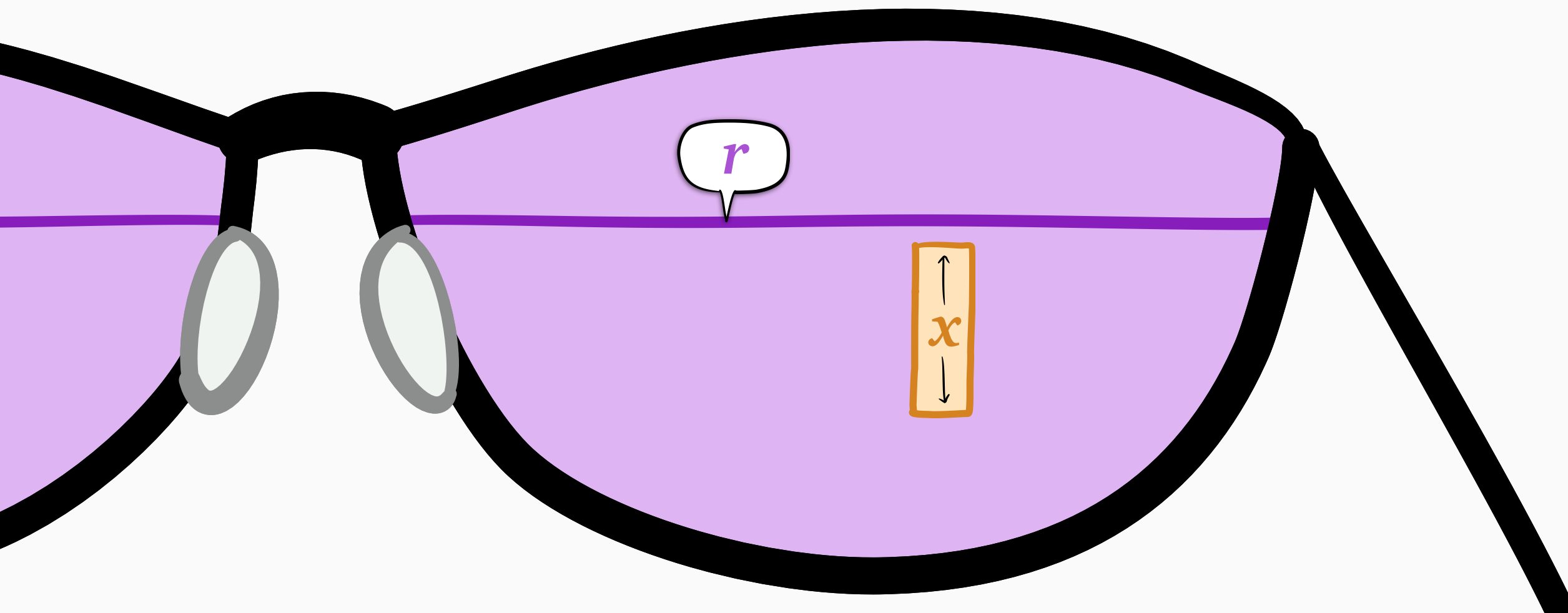
$w_x(r)$ = r -work of *single job* of rem. size x = $\begin{cases} 0 & \text{if } r < x \end{cases}$



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

$$w_x(r) = r\text{-work of single job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$



Defining one job's r -work

$W(r)$ = work relevant to **rank** r

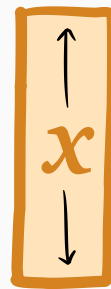
$$w_x(r) = r\text{-work of single job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$



Defining one job's r -work

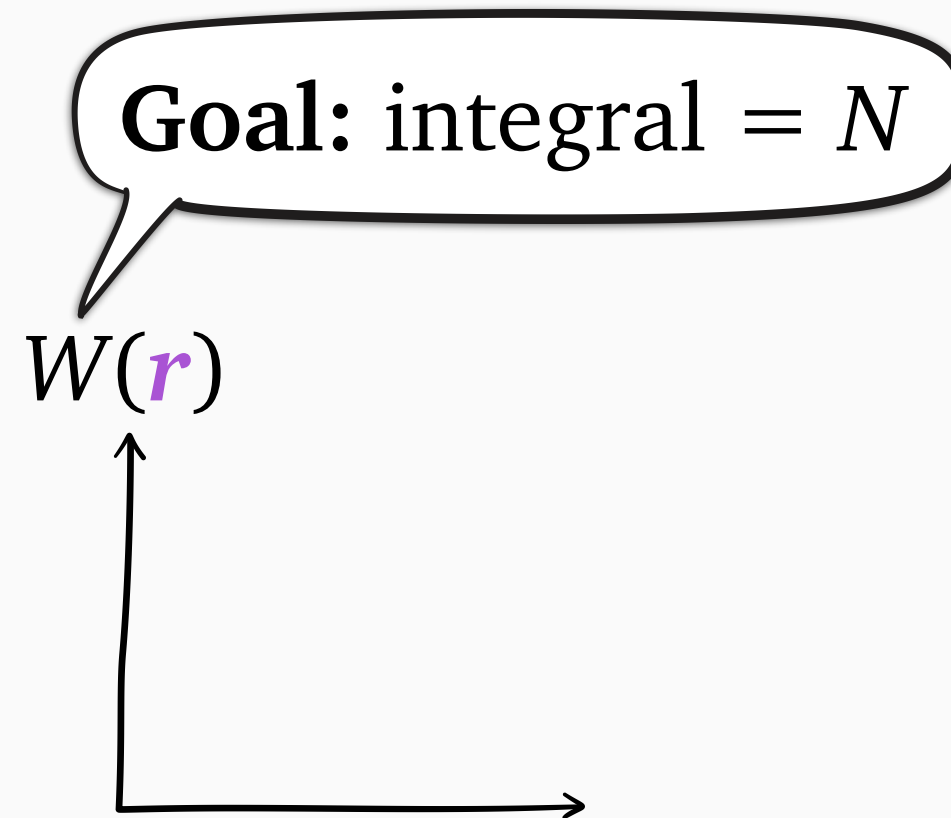
$W(r)$ = work relevant to **rank** r
= total r -work of all jobs

$w_x(r)$ = r -work of *single job* of rem. size x = $\begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$

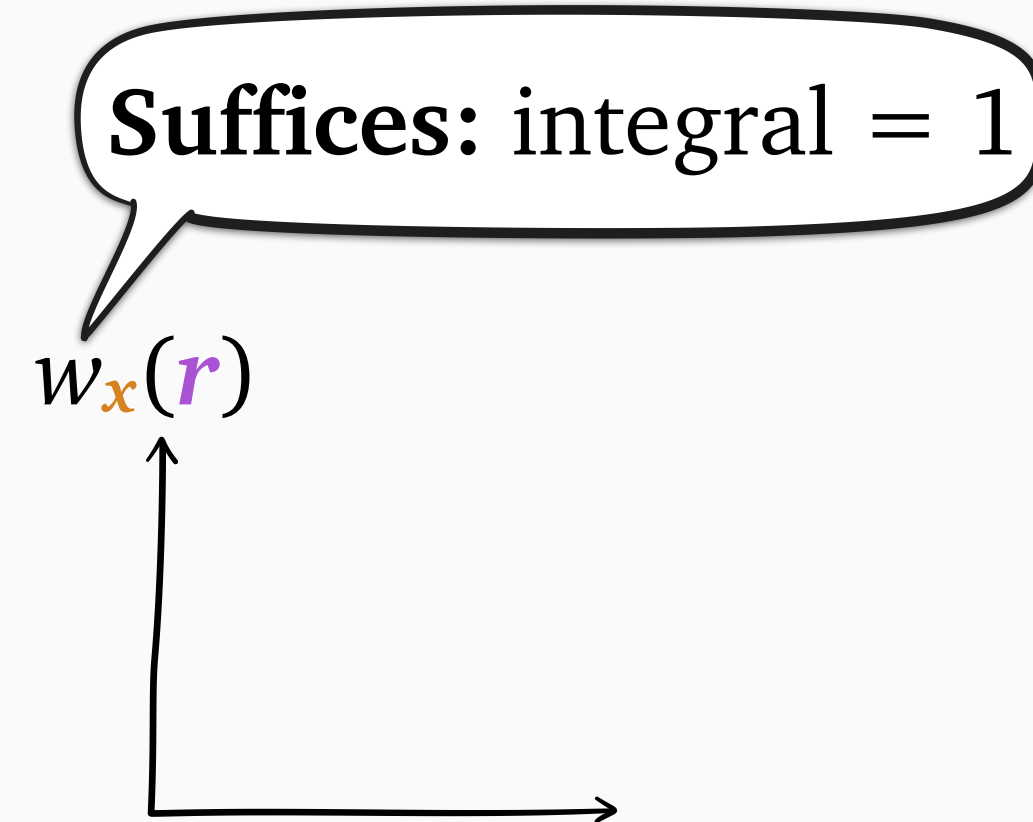
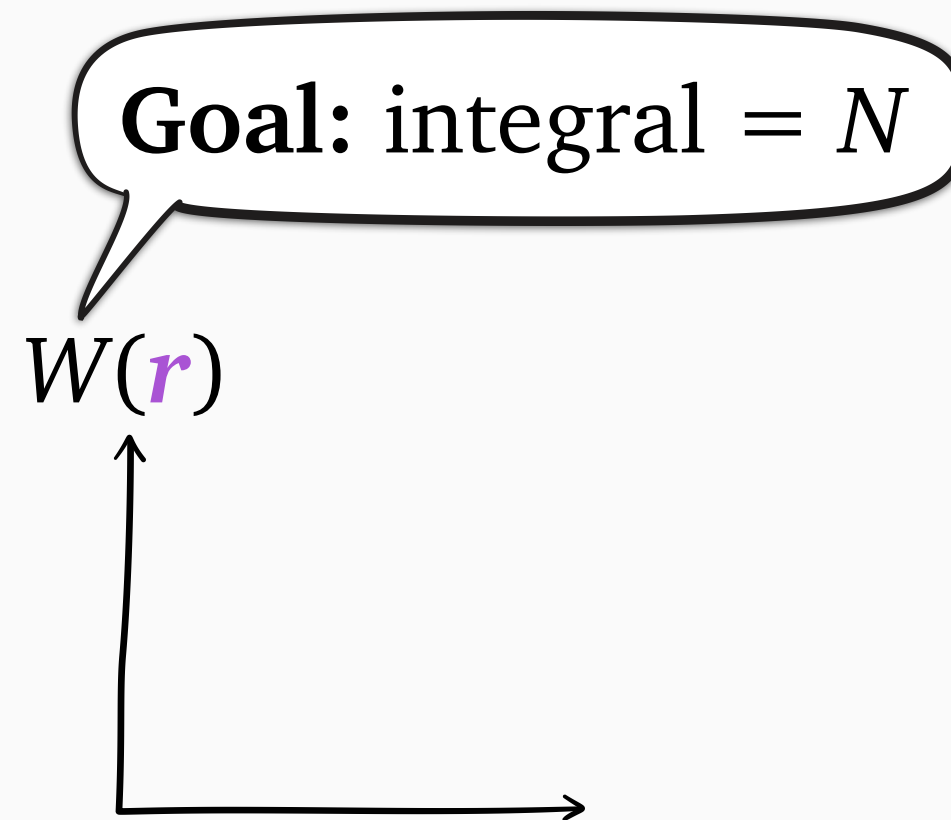


From *r*-work to number of jobs

From r -work to number of jobs

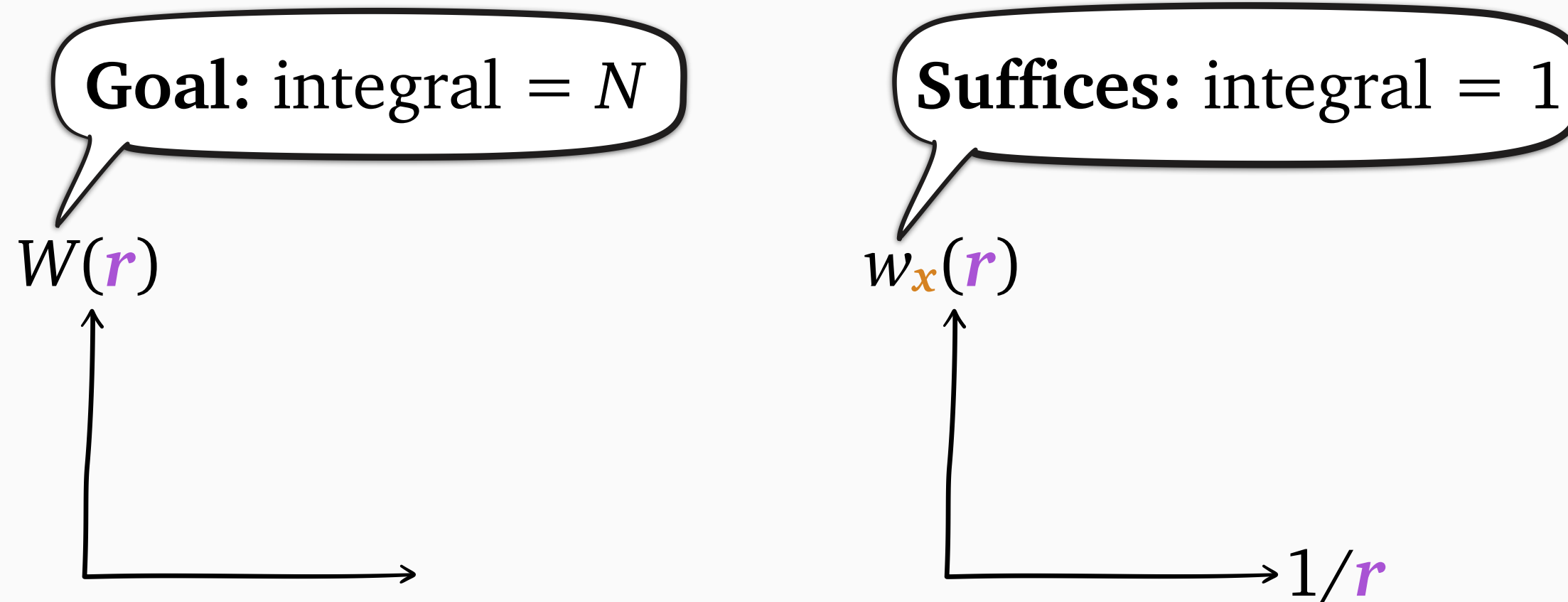


From r -work to number of jobs



$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

From r -work to number of jobs



$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

From r -work to number of jobs

Goal: integral = N

$W(r)$



Suffices: integral = 1

$w_x(r)$



$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

From r -work to number of jobs

Goal: integral = N

$W(r)$



Suffices: integral = 1

$w_x(r)$

x

1

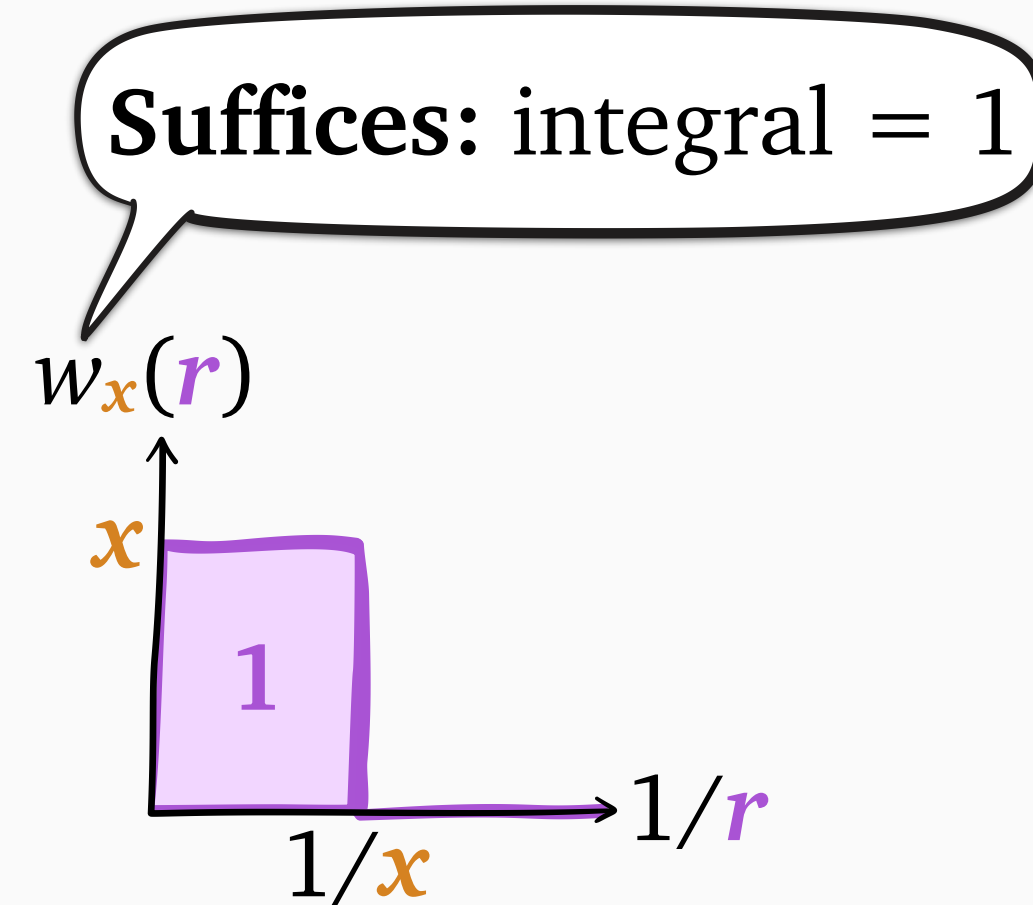
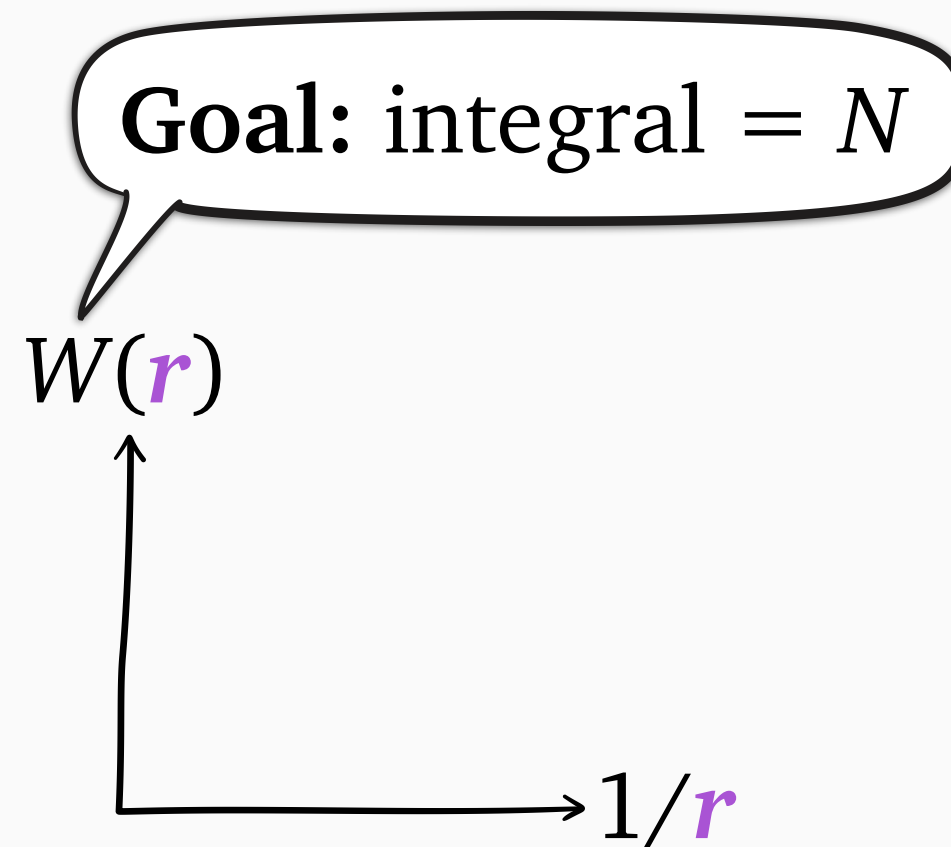
$1/x$

$1/r$



$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

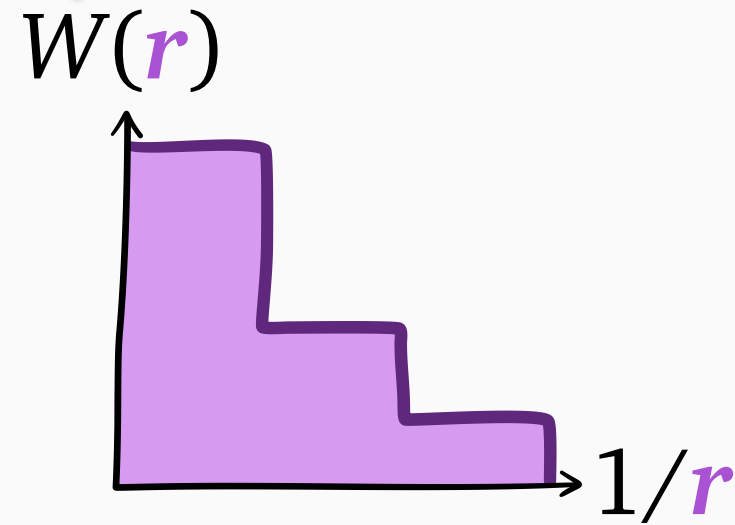
From r -work to number of jobs



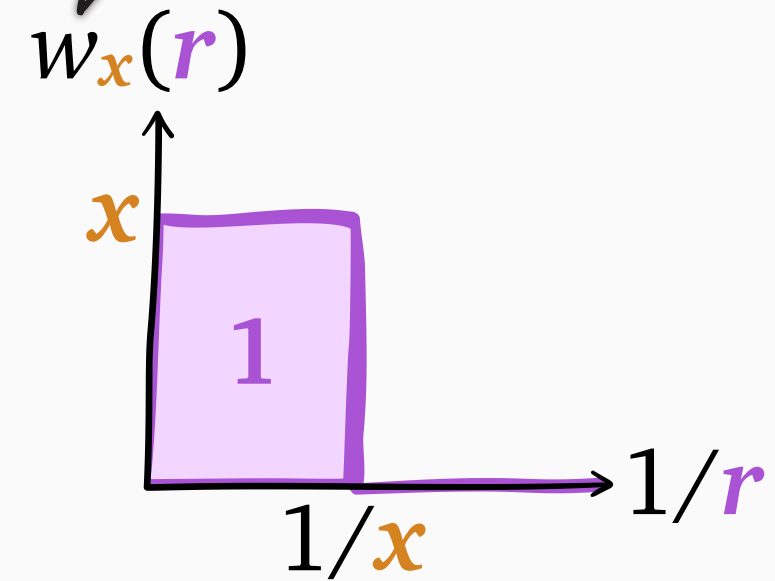
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

From r -work to number of jobs

Goal: integral = N



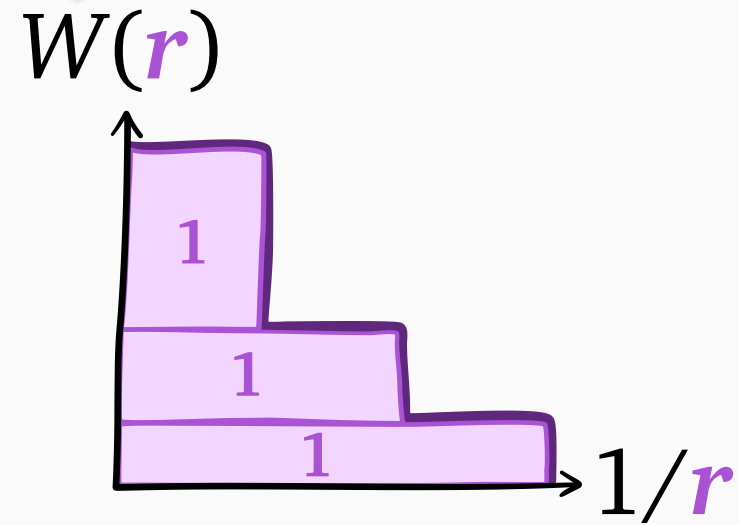
Suffices: integral = 1



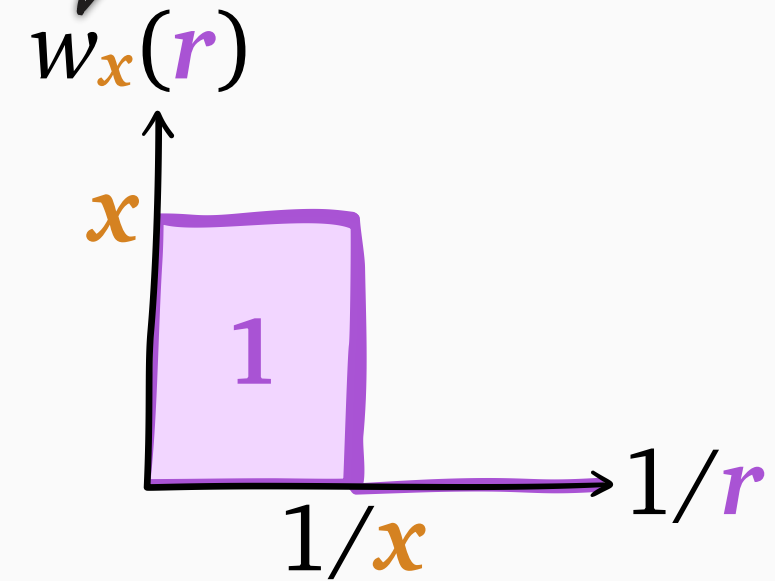
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

From r -work to number of jobs

Goal: integral = N



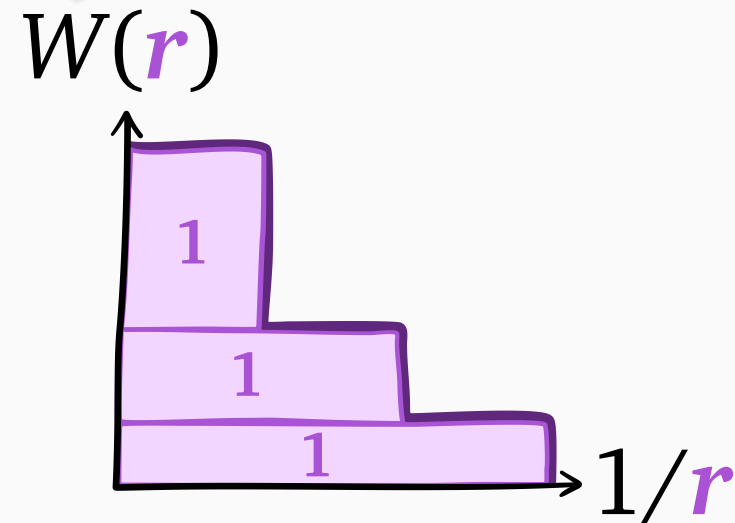
Suffices: integral = 1



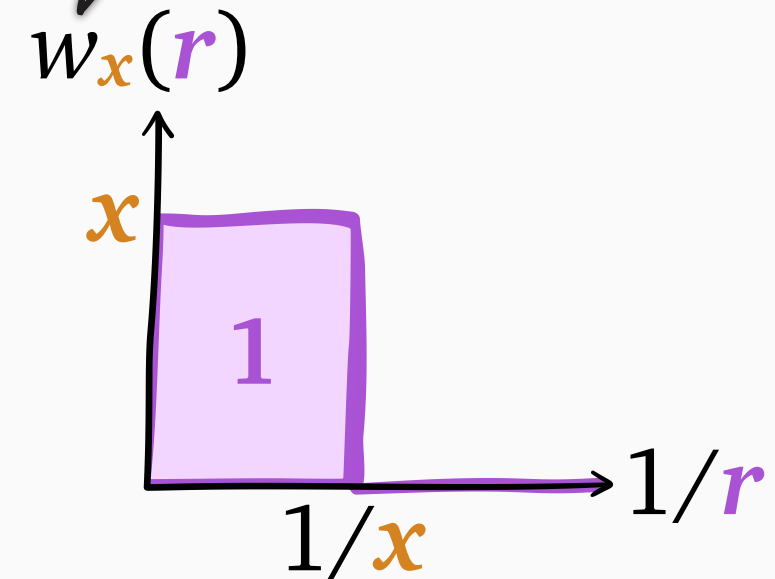
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

From r -work to number of jobs

Goal: integral = N



Suffices: integral = 1

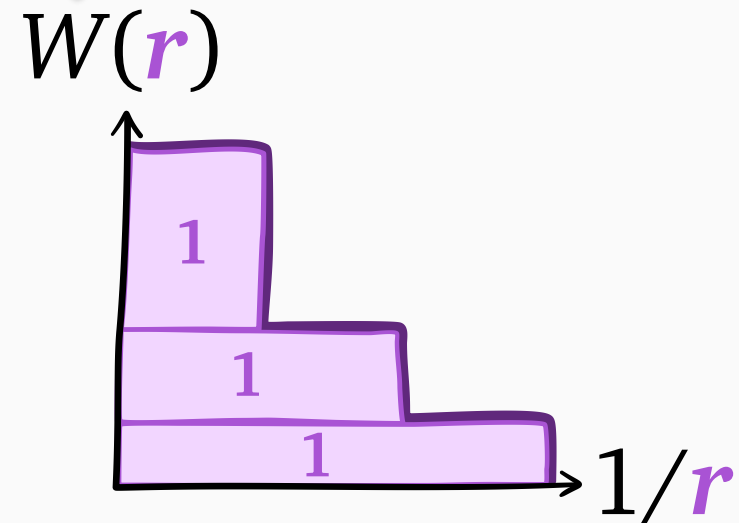


Theorem:

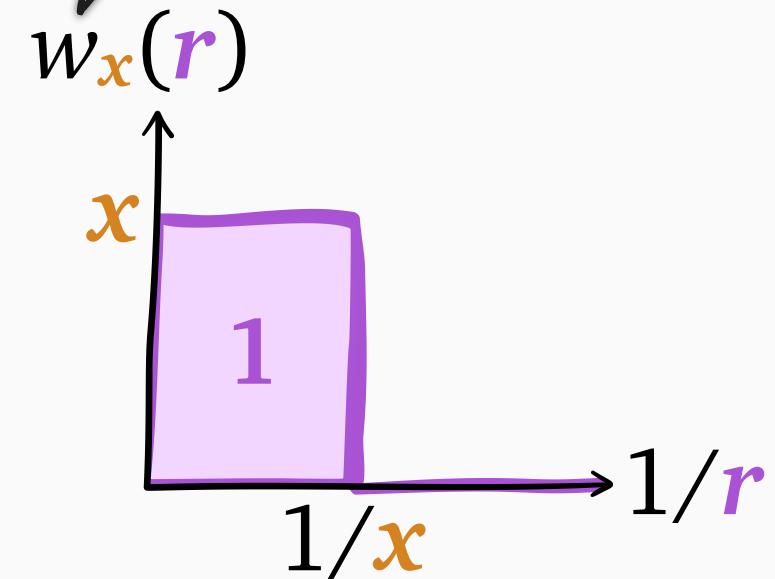
$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

From r -work to number of jobs

Goal: integral = N



Suffices: integral = 1



Theorem:

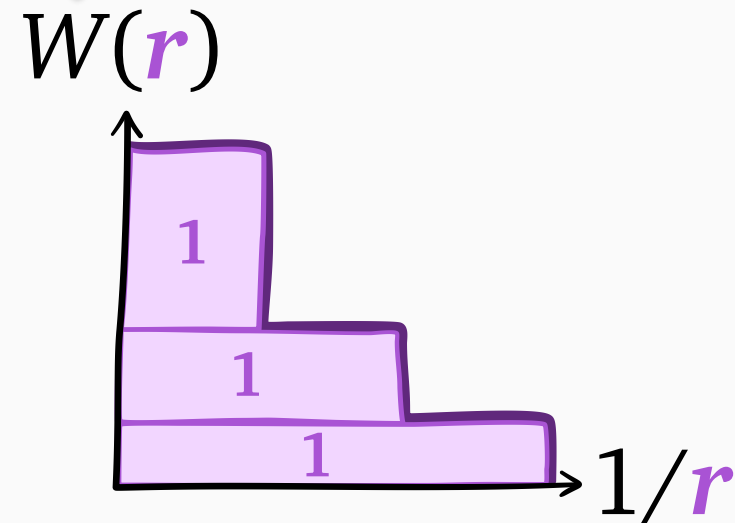


WINE

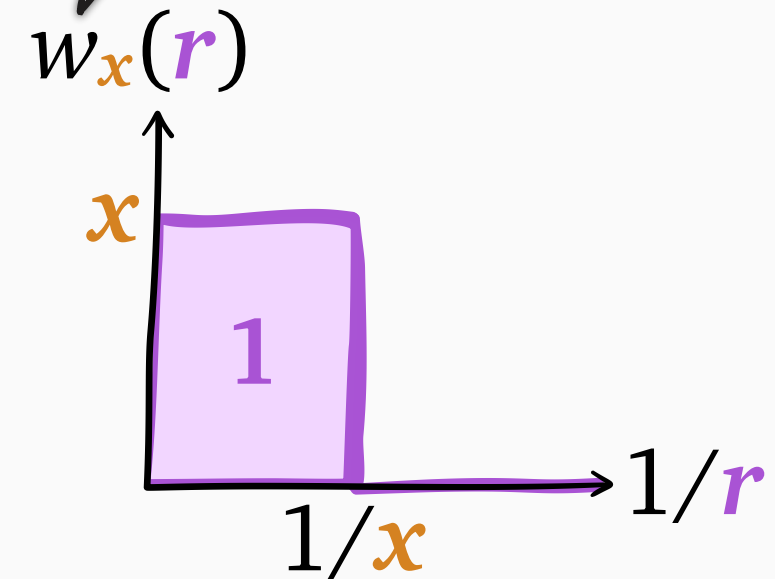
$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

From r -work to number of jobs

Goal: integral = N



Suffices: integral = 1



Theorem:



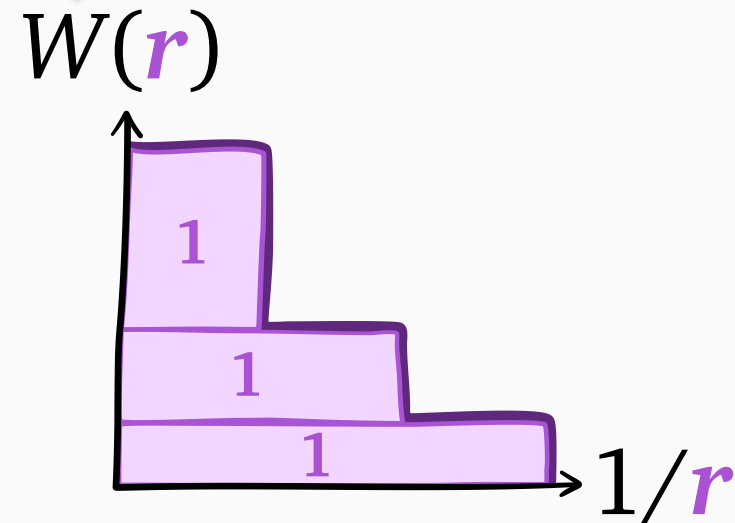
WINE

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

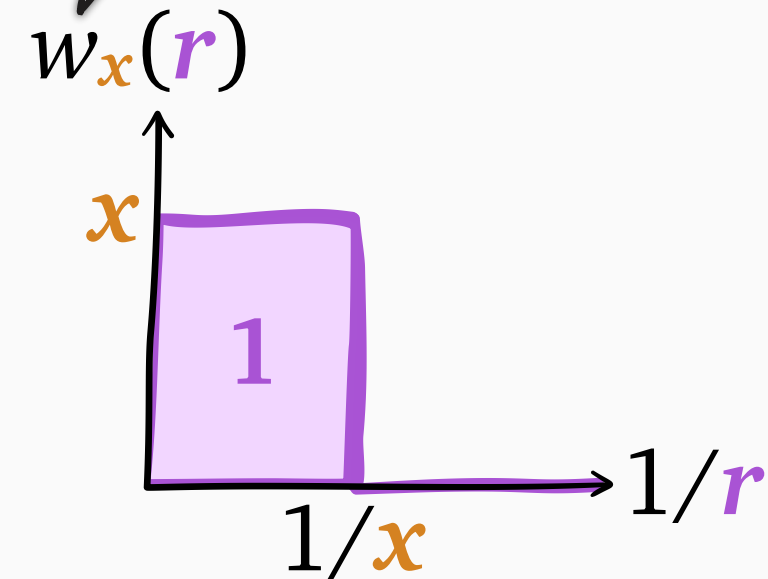
uses rank = rem. size

From r -work to number of jobs

Goal: integral = N



Suffices: integral = 1



Theorem: under *any* policy,




WINE

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

uses **rank** = rem. size


How does WINE help?

 WINE

Theorem:

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

How does WINE help?


 WINE

Theorem:

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

How to minimize $W(r)$?

How does WINE help?

 WINE


Theorem:

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

How to minimize $W(r)$?

*Prioritize jobs with **rank** $\leq r$*

How does WINE help?

 **Theorem:**

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$


How to minimize $W(r)$?

Prioritize jobs with $\text{rank} \leq r$

To do for all r :

*always serve job of minimum **rank***

How does **WINE** help?

 **Theorem:**

$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

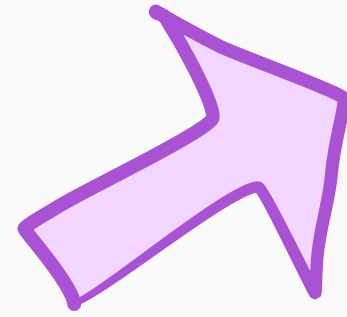
How to minimize $W(r)$?

Prioritize jobs with **rank** $\leq r$

To do for all r :
*always serve job of minimum **rank***

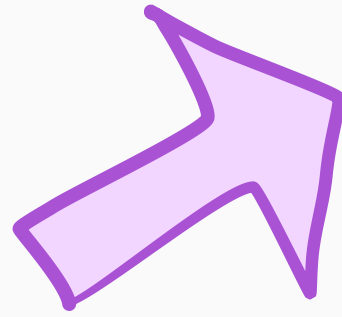


under Poisson
arrivals



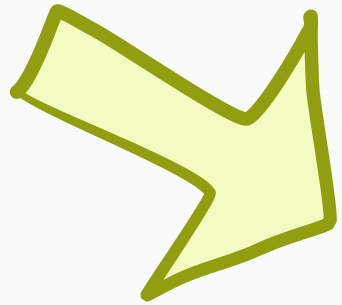
SRPT-flavored WINE:

$$N = \int_0^{\infty} \frac{W(\mathbf{r})}{\mathbf{r}^2} d\mathbf{r}$$



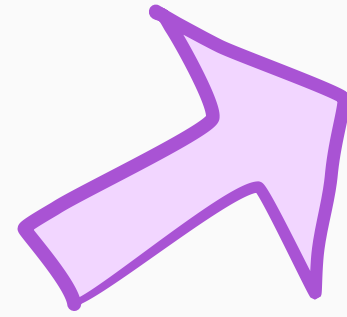
SRPT-flavored WINE:

$$N = \int_0^\infty \frac{W(\mathbf{r})}{\mathbf{r}^2} d\mathbf{r}$$



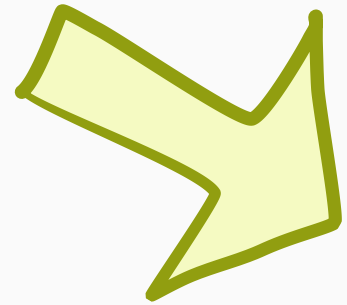
Gittins-flavored WINE:

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



SRPT-flavored WINE:

$$N = \int_0^{\infty} \frac{W(\mathbf{r})}{\mathbf{r}^2} d\mathbf{r}$$

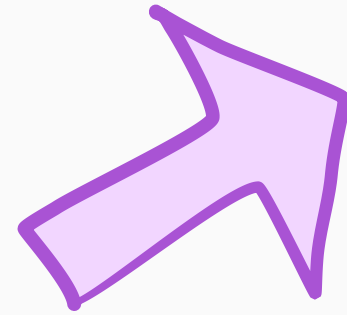
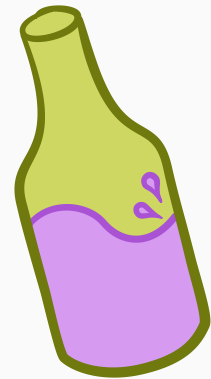


Gittins-flavored WINE:

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

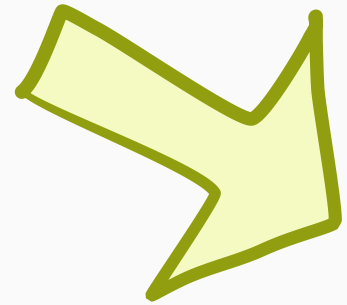
Lemma: using **Gittins**-flavored \mathbf{r} -work,

$$1 = \int_0^{\infty} \frac{\mathbf{E}[\text{one job's } \mathbf{r}\text{-work} \mid \text{job's state}]}{\mathbf{r}^2} d\mathbf{r}$$



SRPT-flavored WINE:

$$N = \int_0^{\infty} \frac{W(\mathbf{r})}{\mathbf{r}^2} d\mathbf{r}$$



Gittins-flavored WINE:

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

Related to *achievable
region method*

[Bertsimas & Niño-Mora]

Lemma: using **Gittins**-flavored \mathbf{r} -work,

$$1 = \int_0^{\infty} \frac{\mathbf{E}[\text{one job's } \mathbf{r}\text{-work} \mid \text{job's state}]}{\mathbf{r}^2} d\mathbf{r}$$

Robustness of Gittins

$$\mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(\textcolor{violet}{r})]}{\textcolor{violet}{r}^2} d\textcolor{violet}{r}$$

Robustness of **Gittins**

$$\mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(\textcolor{violet}{r})]}{\textcolor{violet}{r}^2} d\textcolor{violet}{r}$$

Corollary: if **rank** function is within c factor of **Gittins**'s, then $\mathbf{E}[N]$ is within c^2 of optimal

Robustness of **Gittins**

$$\mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(\textcolor{violet}{r})]}{\textcolor{violet}{r}^2} d\textcolor{violet}{r}$$

Corollary: if **rank** function is within c factor of **Gittins**'s, then $\mathbf{E}[N]$ is within c^2 of optimal

Proof: change of variables in integral

Robustness of **Gittins**

$$\mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(\textcolor{violet}{r})]}{\textcolor{violet}{r}^2} d\textcolor{violet}{r}$$

robustness to noisy
job size predictions
[Scully et al., 2022]

Corollary: if **rank** function is within c factor
of **Gittins**'s, then $\mathbf{E}[N]$ is within c^2 of optimal

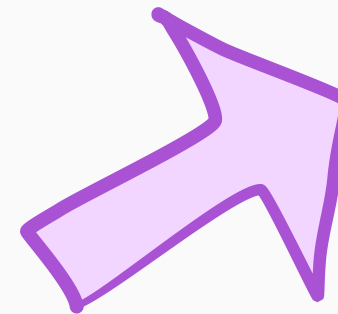
Proof: change of variables in integral

Robustness of Gittins

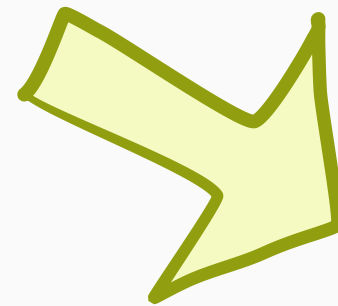
$$\mathbf{E}[N] = \int_0^\infty \frac{\mathbf{E}[W(r)]}{r^2} dr$$

Corollary: if **rank** function is within c factor of **Gittins**'s, then $\mathbf{E}[N]$ is within c^2 of optimal

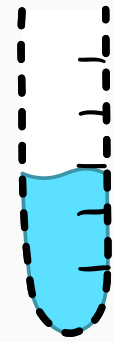
Proof: change of variables in integral



robustness to noisy
job size predictions
[Scully et al., 2022]

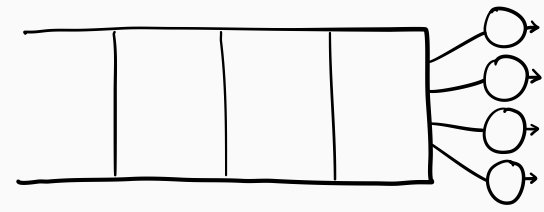


can substitute n samples
for true distribution S
[Ramakrishna et al., 2025]



Part I

Handling job size uncertainty



Part II

Analyzing multiserver scheduling



Part III

Optimizing tail metrics

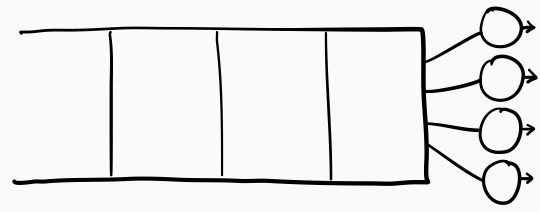


Part I

Handling job size uncertainty

Queueing for TCS

Use WINE to analyze Gittins
with arbitrary release dates?



Part II

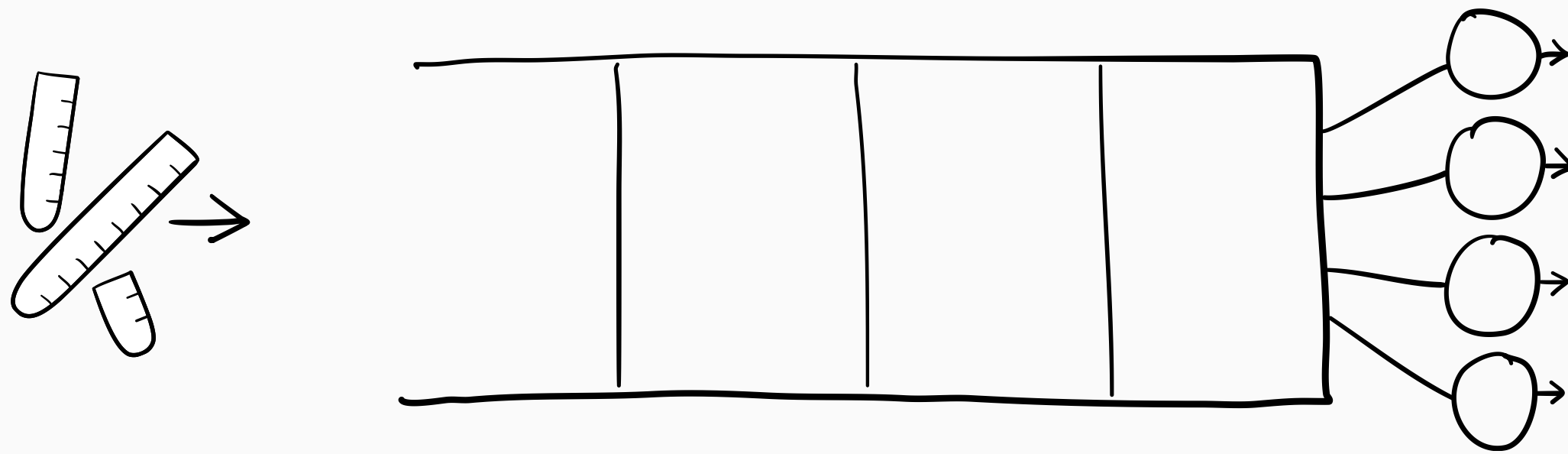
Analyzing multiserver scheduling



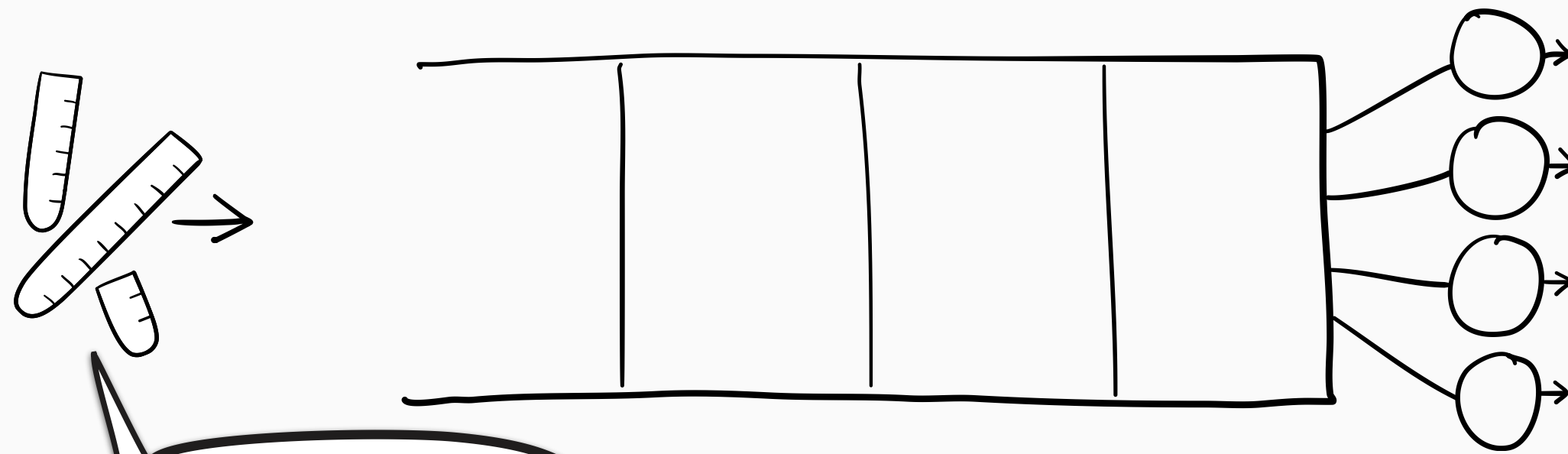
Part III

Optimizing tail metrics

Scheduling in the **M/G/k**



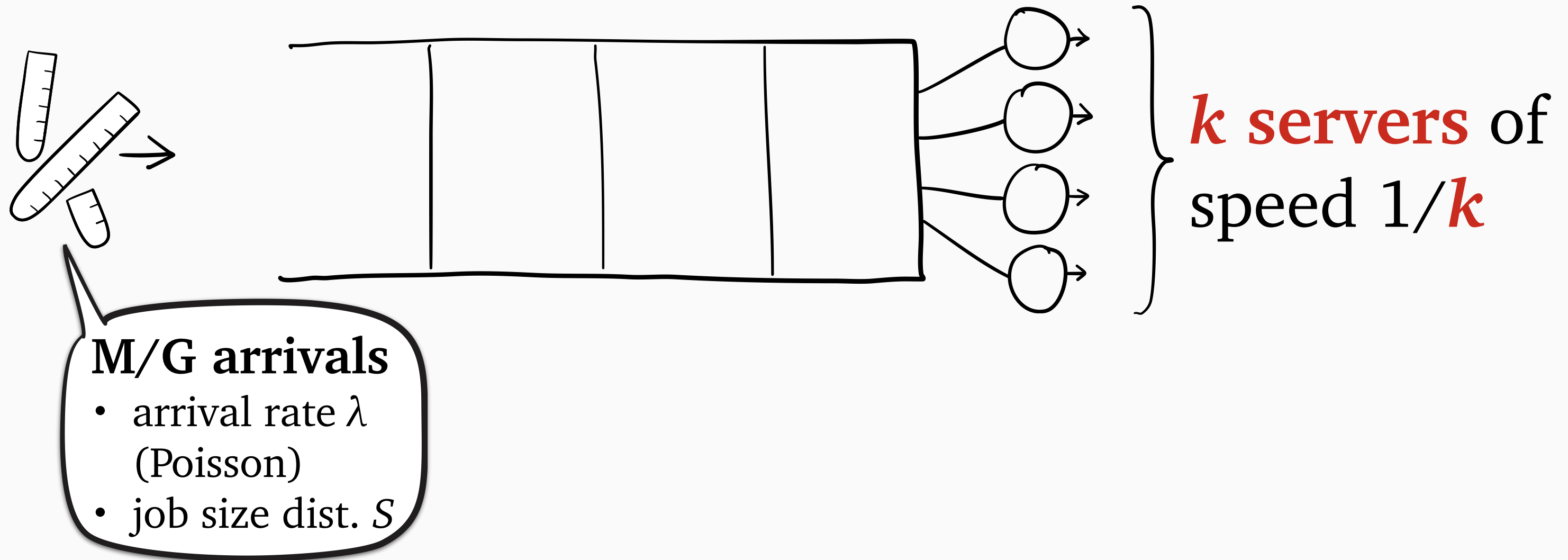
Scheduling in the **M/G/k**



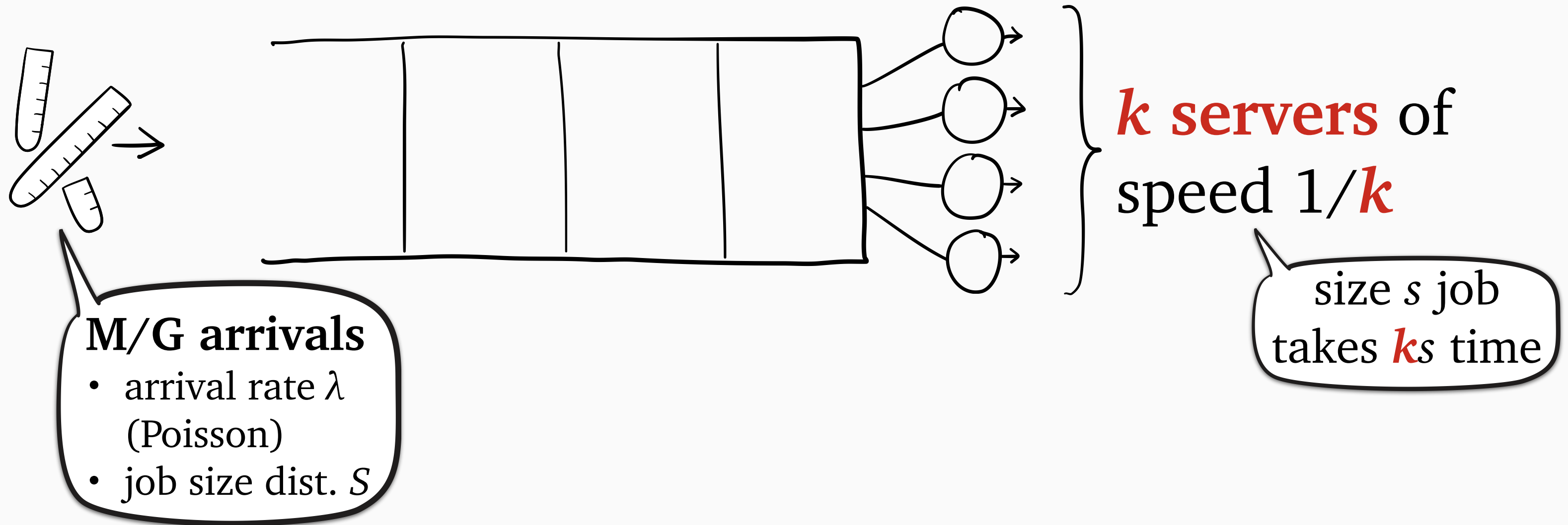
M/G arrivals

- arrival rate λ
(Poisson)
- job size dist. S

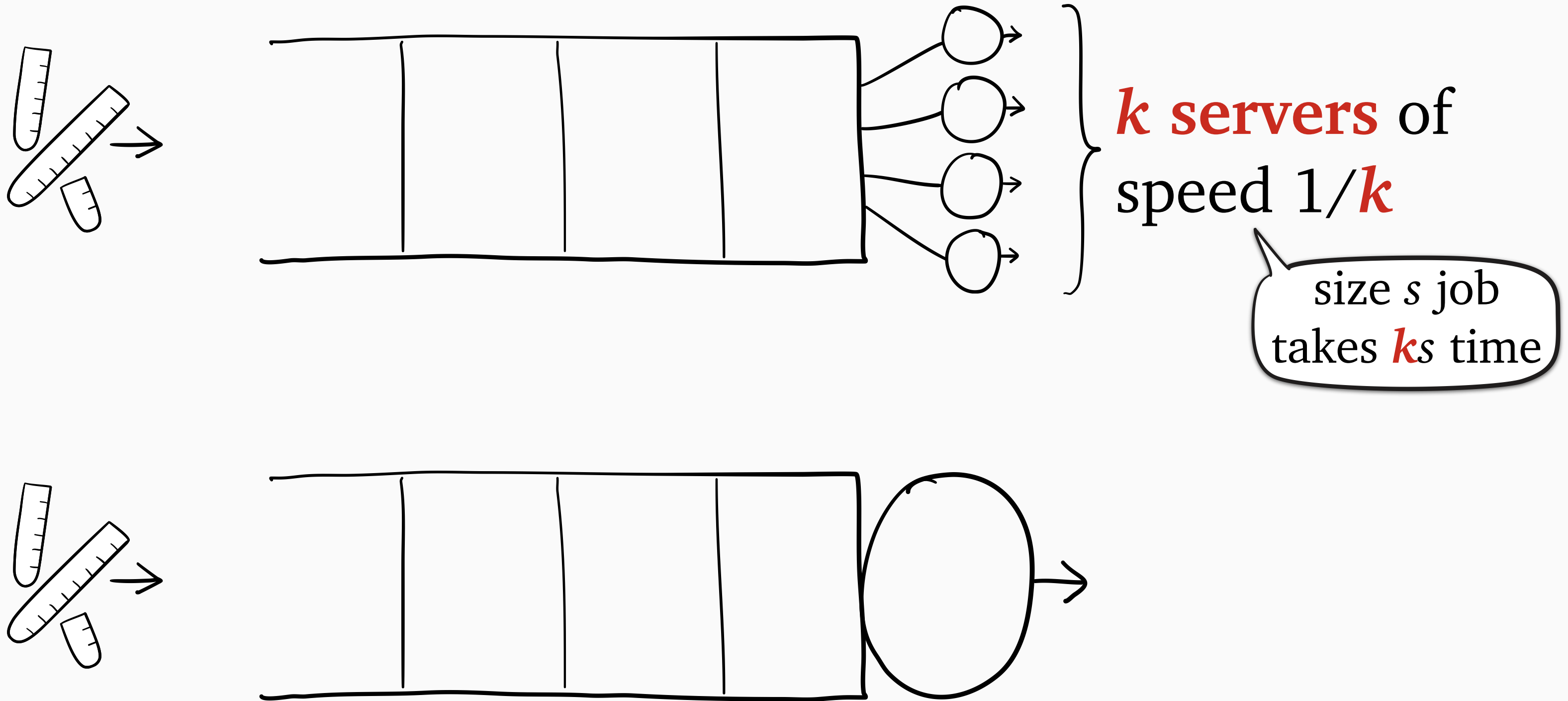
Scheduling in the **M/G/k**



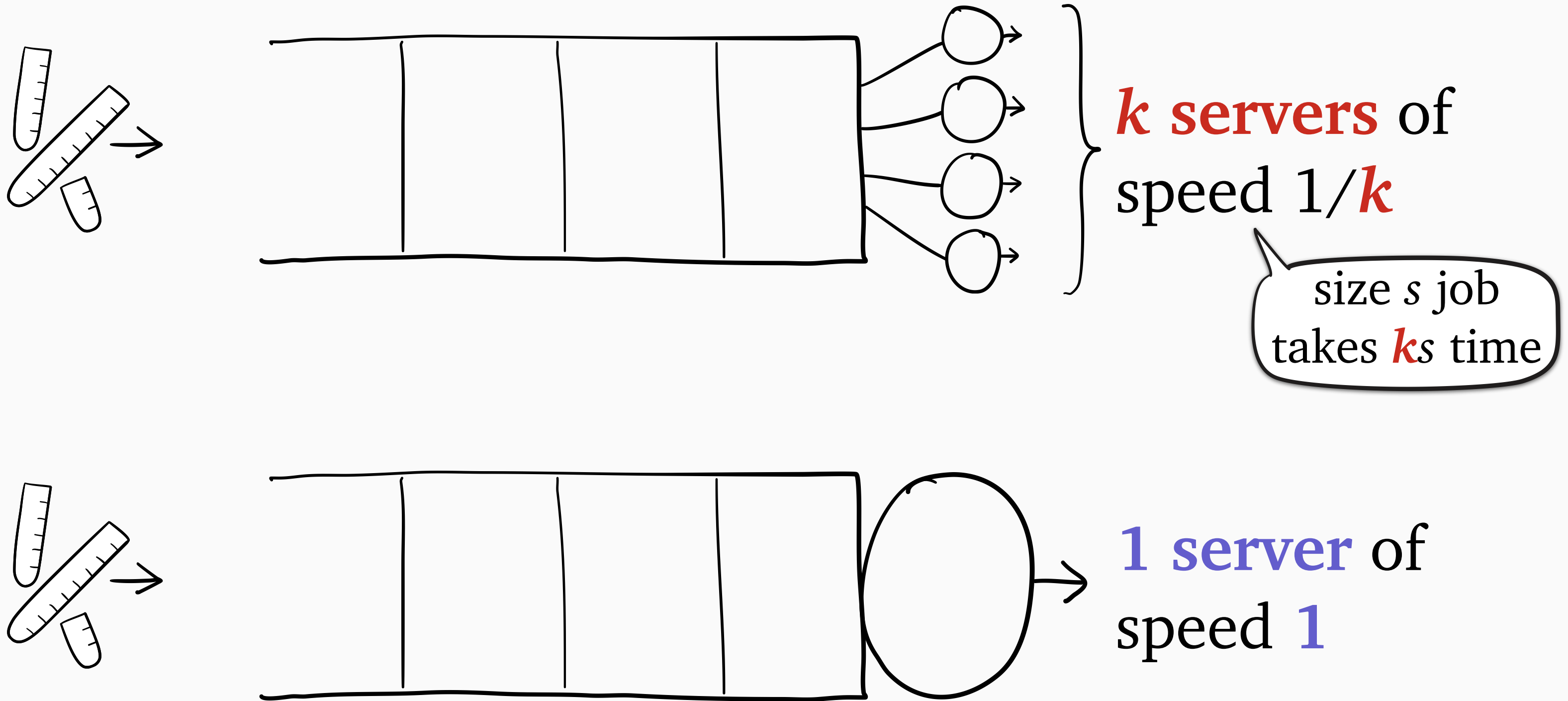
Scheduling in the **M/G/k**



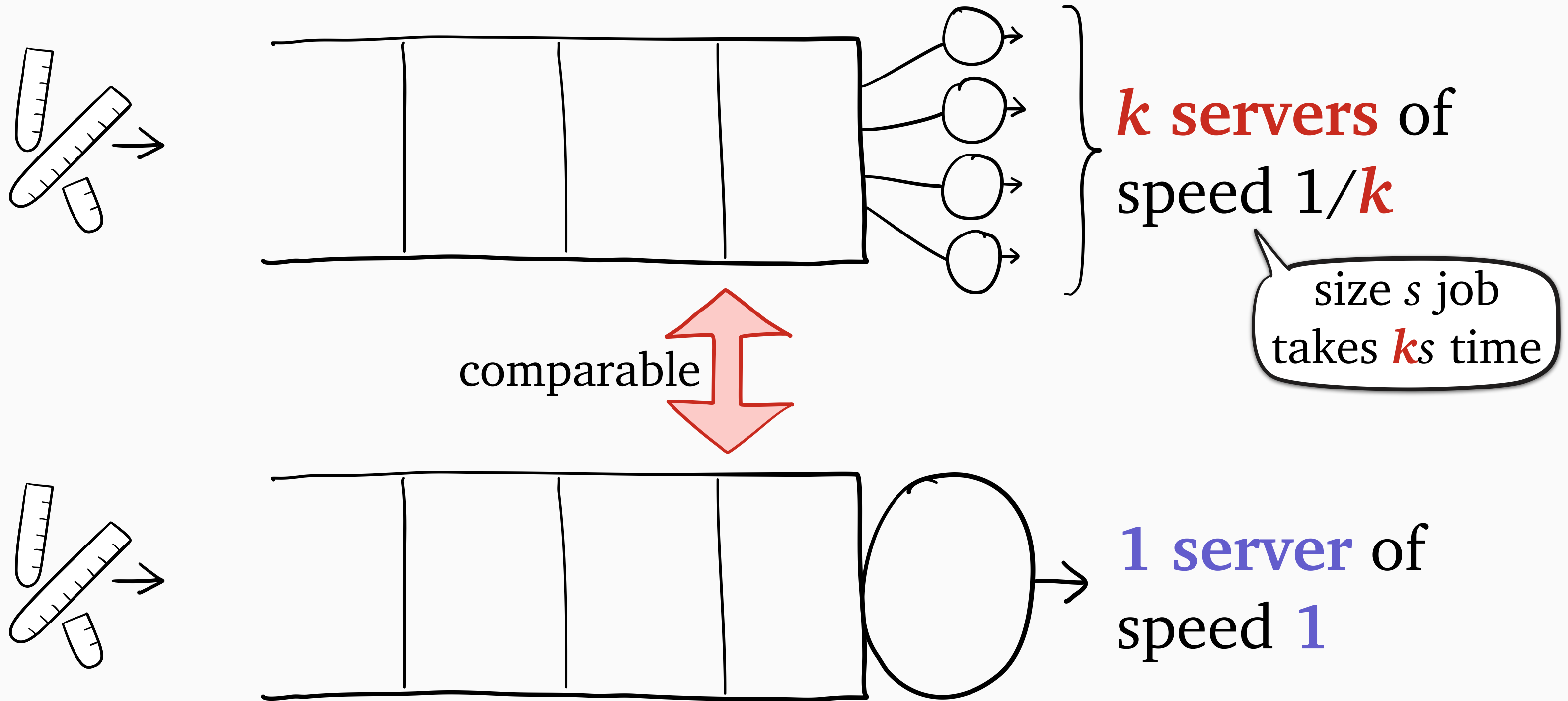
Scheduling in the **M/G/k**



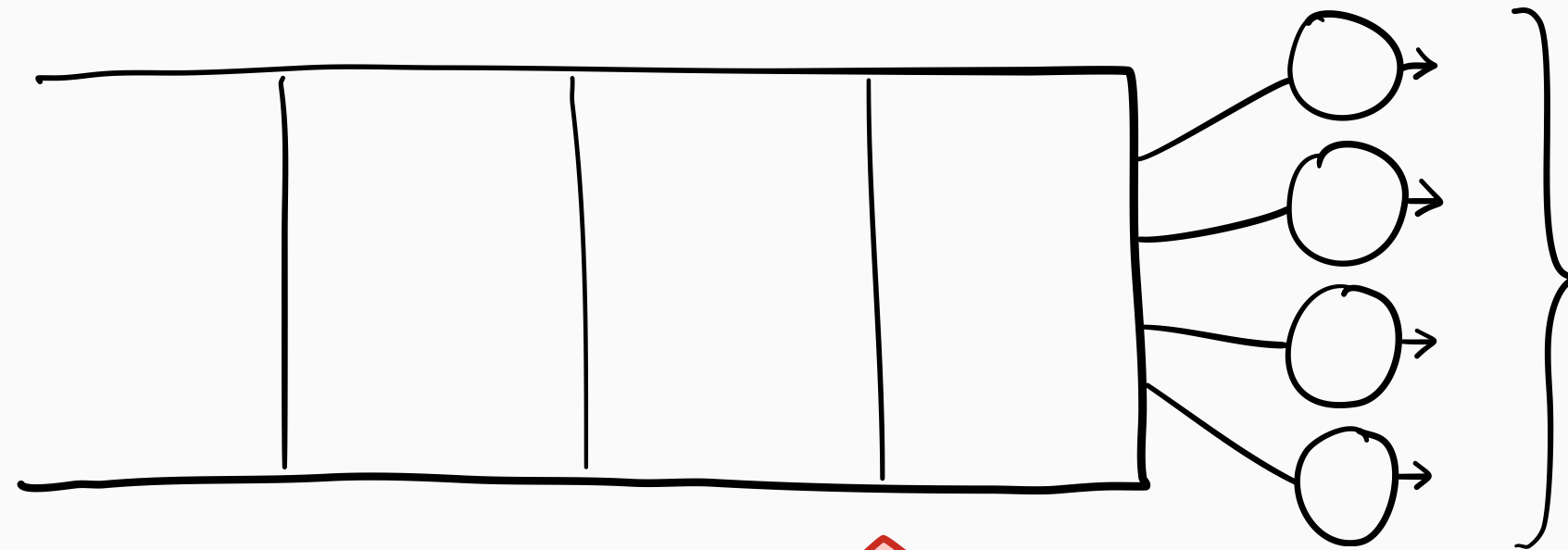
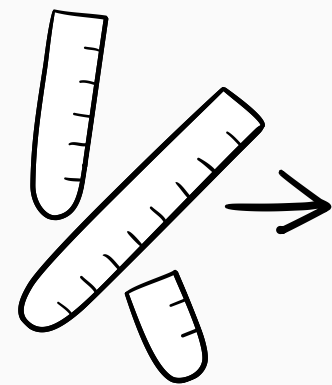
Scheduling in the **M/G/k**



Scheduling in the **M/G/k**



Scheduling in the **M/G/k**

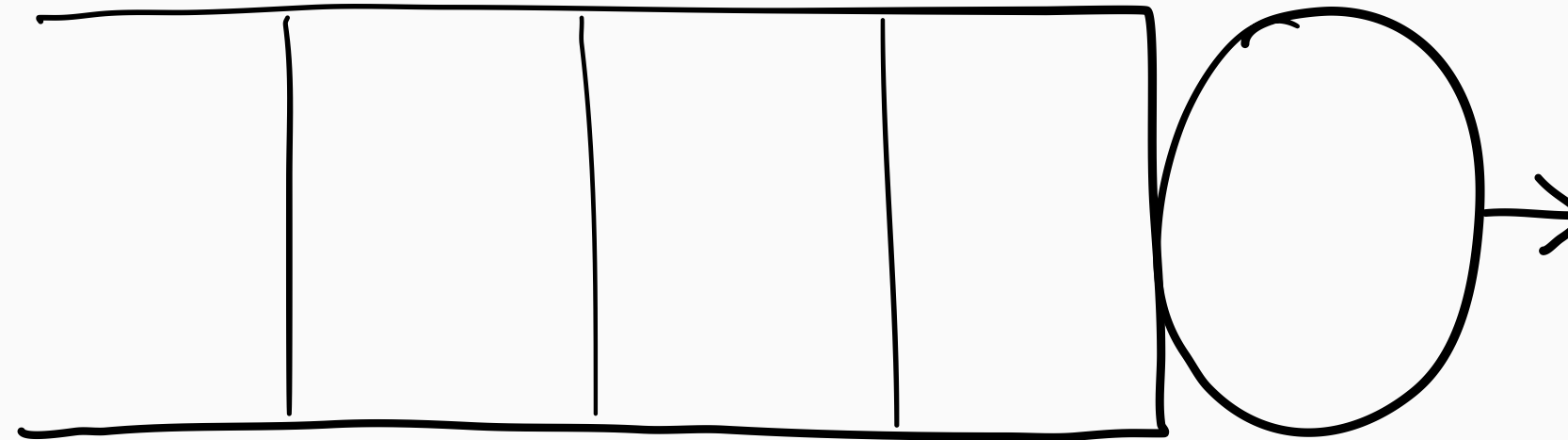
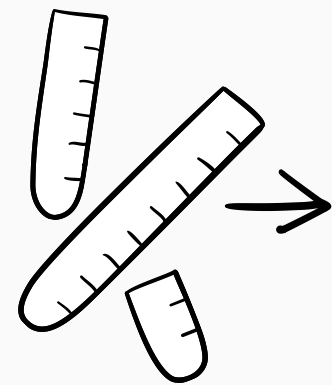
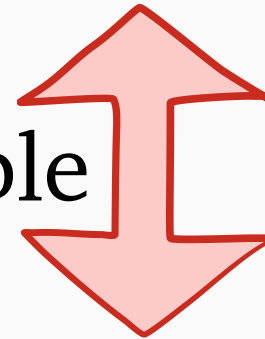


k servers of
speed $1/k$

 bin-packing
aspect

size s job
takes ks time

comparable



1 server of
speed **1**

? Q: Are **SRPT** and **Gittins**
still good in the **M/G/k**?

? Q: Are **SRPT** and **Gittins** still good in the **M/G/k**?

✓ A: Yes, especially in the *heavy traffic limit*

$$\rho = \lambda \mathbf{E}[S] \rightarrow 1$$

? Q: Are **SRPT** and **Gittins** still good in the **M/G/k**?

✓ A: Yes, especially in the *heavy traffic limit*

$$\rho = \lambda \mathbf{E}[S] \rightarrow 1$$

Theorem: for **SRPT** and **Gittins**,

$$\mathbf{E}[T_{\mathbf{k}}] \leq \mathbf{E}[T_{\mathbf{1}}] + (k - 1) \cdot O\left(\log \frac{1}{1 - \rho}\right)$$

? Q: Are **SRPT** and **Gittins** still good in the **M/G/k**?

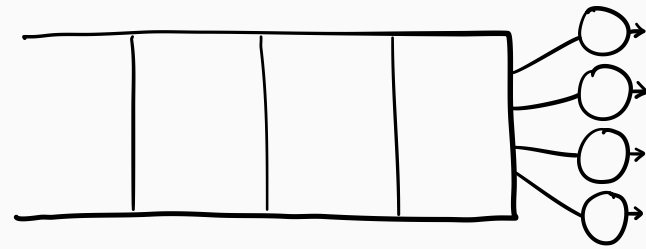
✓ A: Yes, especially in the *heavy traffic limit*

$$\rho = \lambda \mathbf{E}[S] \rightarrow 1$$

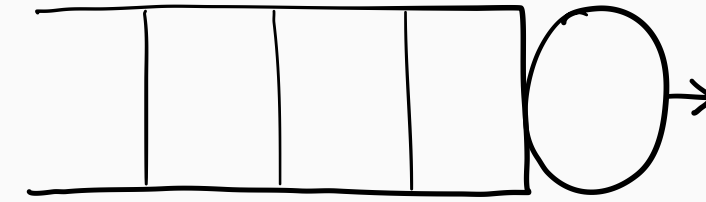
Theorem: for **SRPT** and **Gittins**,

$$\mathbf{E}[T_k] \leq \mathbf{E}[T_1] + \underbrace{(k-1) \cdot O\left(\log \frac{1}{1-\rho}\right)}_{\mathbf{E}[S^2(\log S)^+] < \infty \Rightarrow o(\mathbf{E}[T_1])}$$

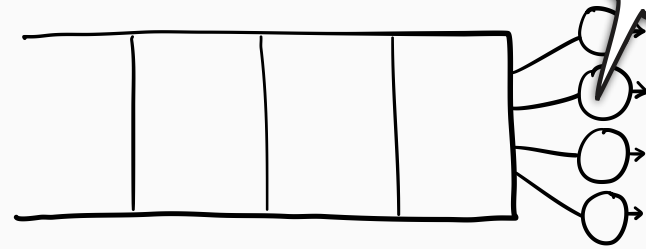
SRPT-*k*



SRPT-1

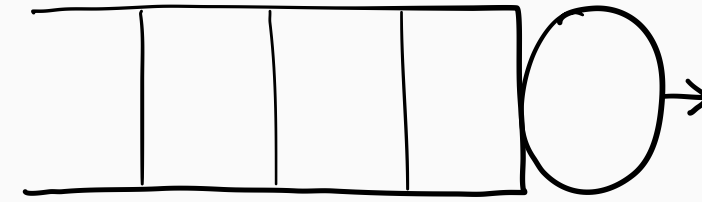


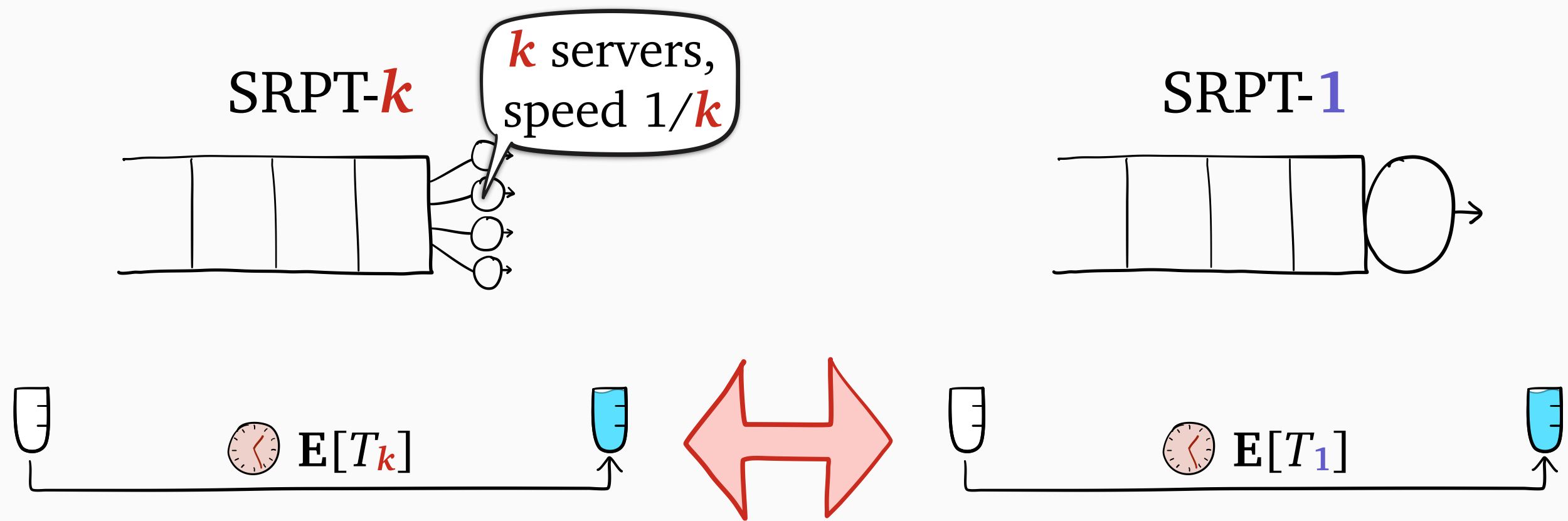
SRPT- k

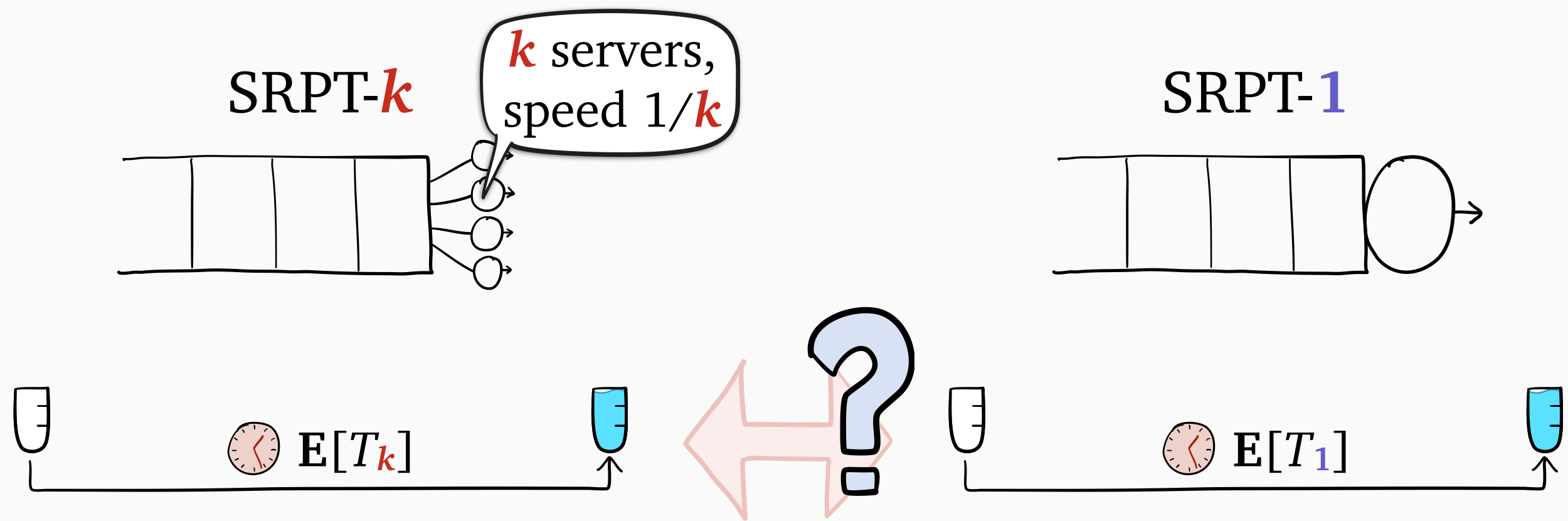


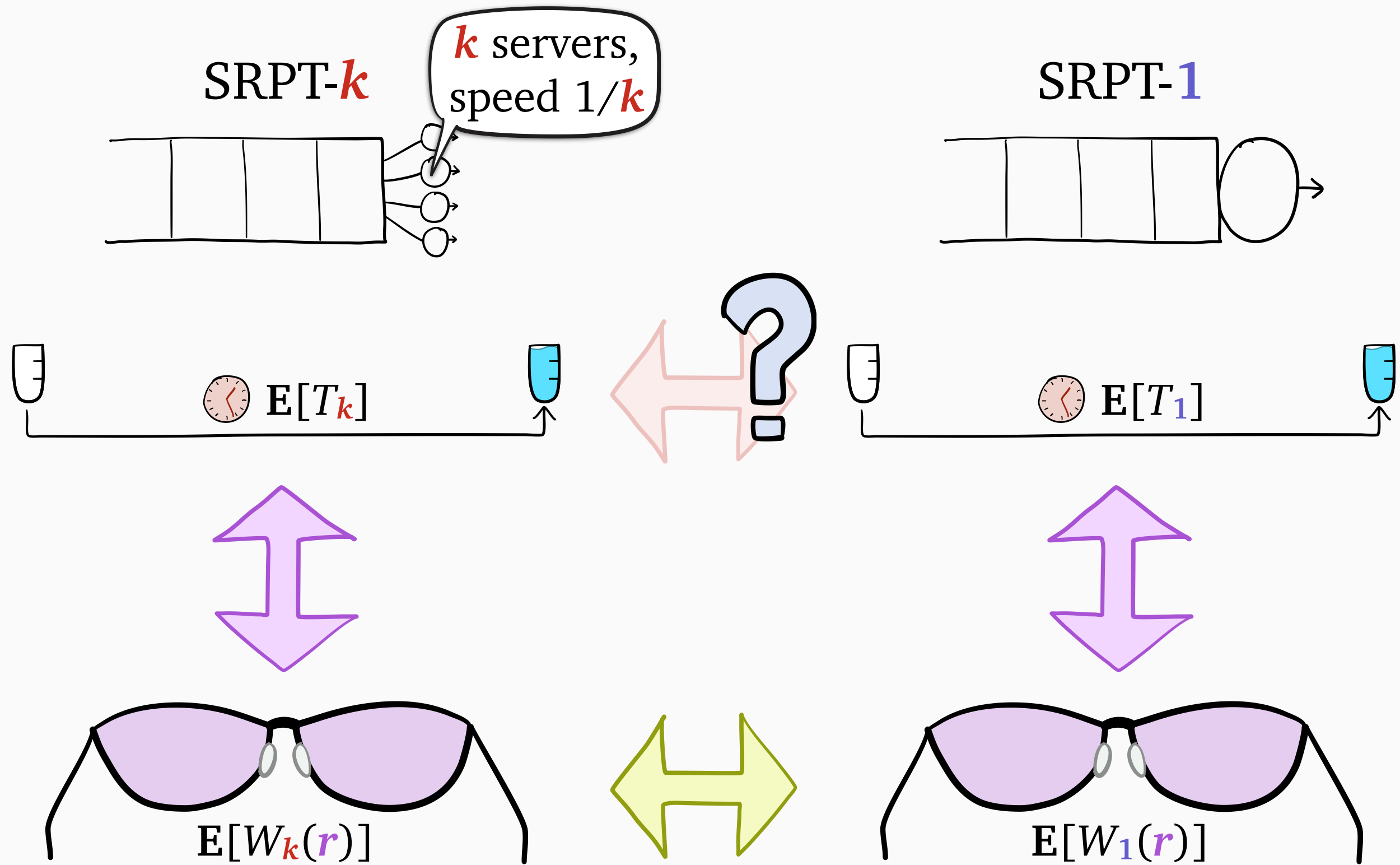
k servers,
speed $1/k$

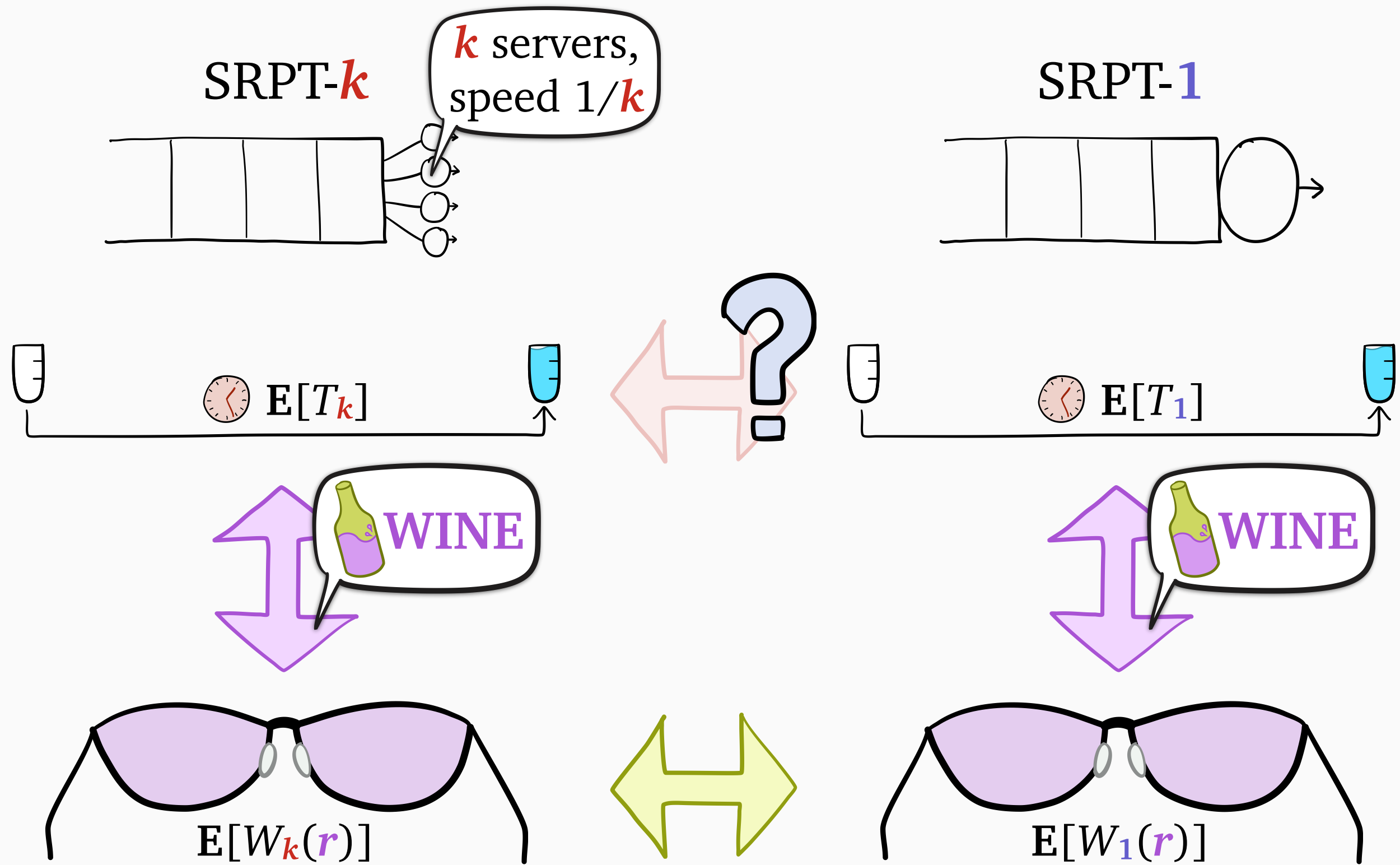
SRPT-1

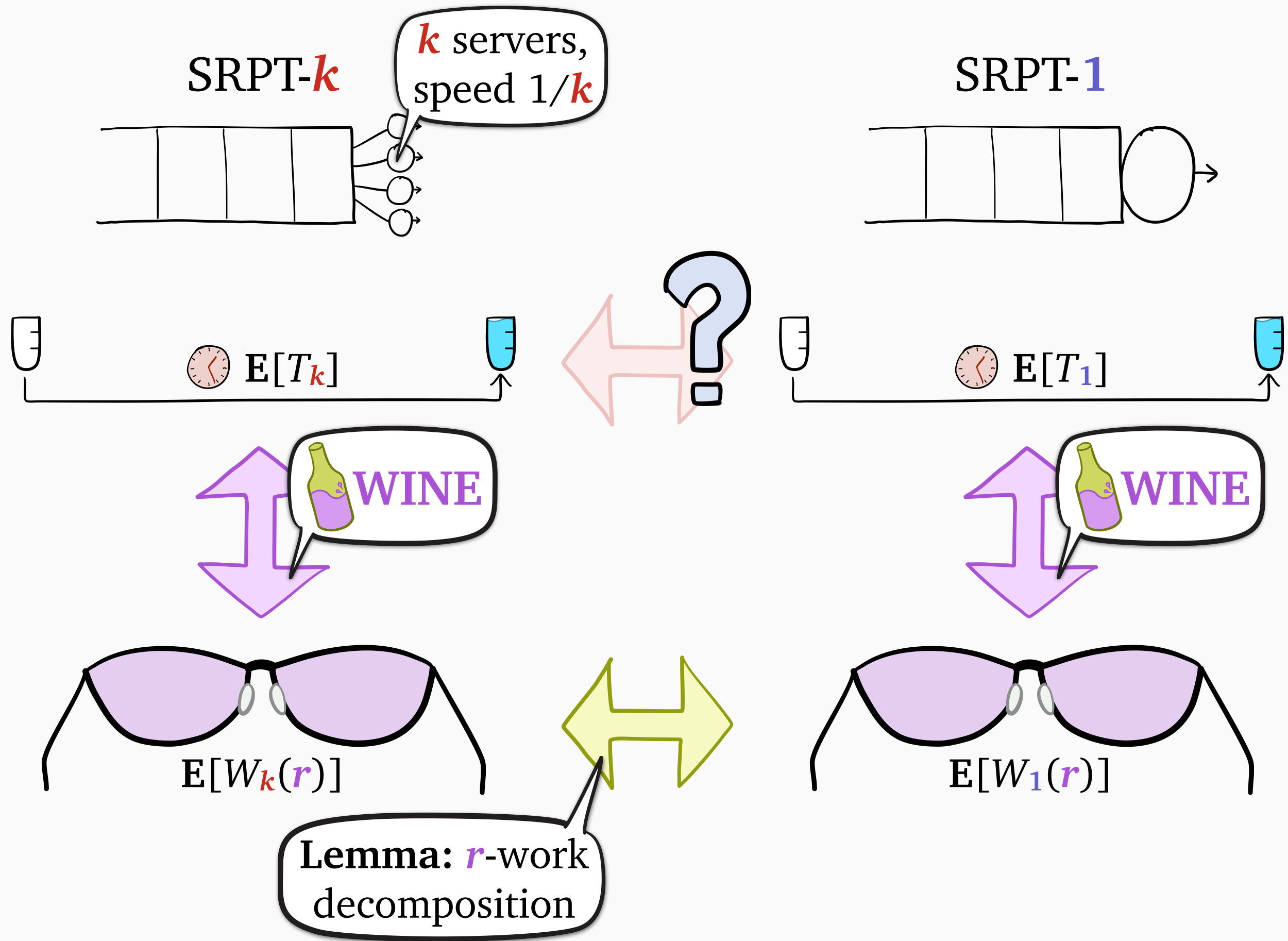


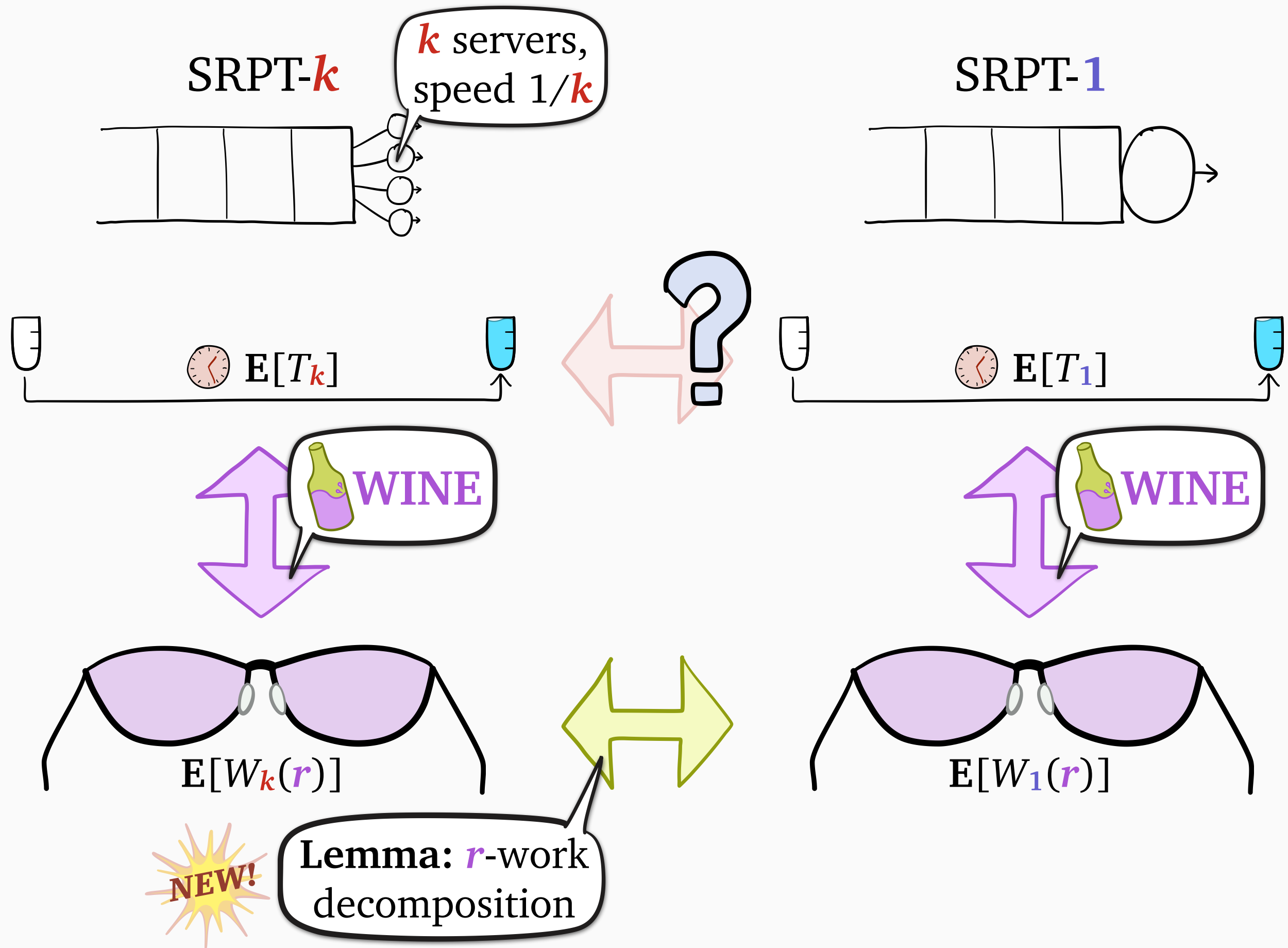


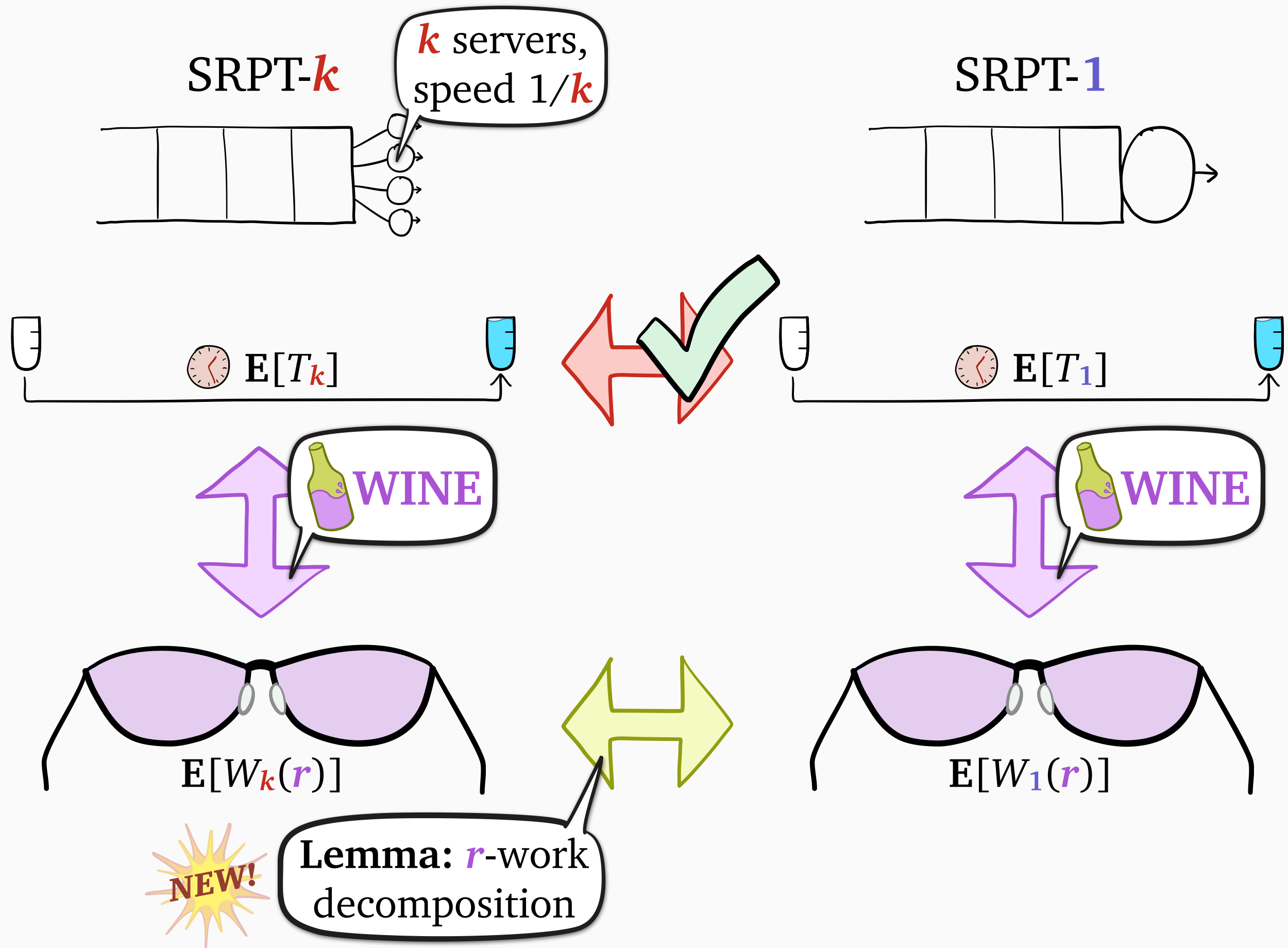






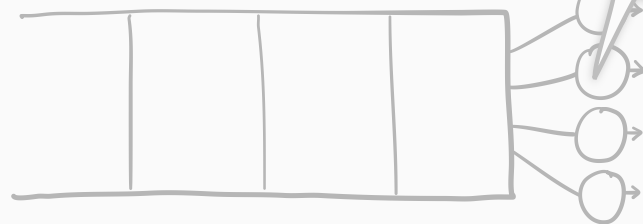




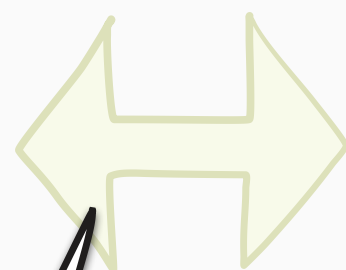
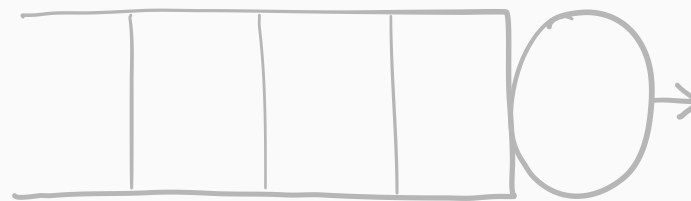


SRPT- k

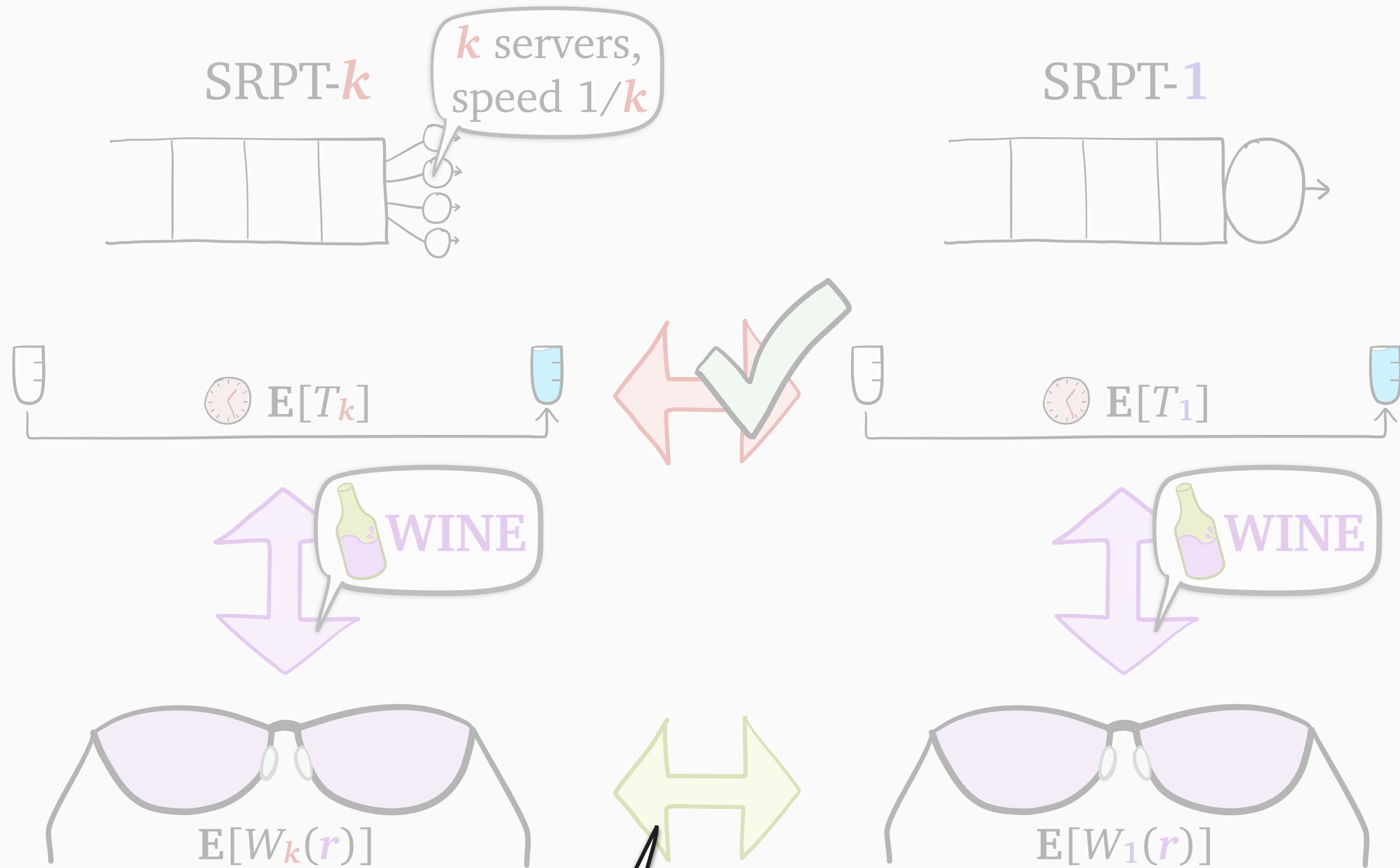
k servers,
speed $1/k$



SRPT-1

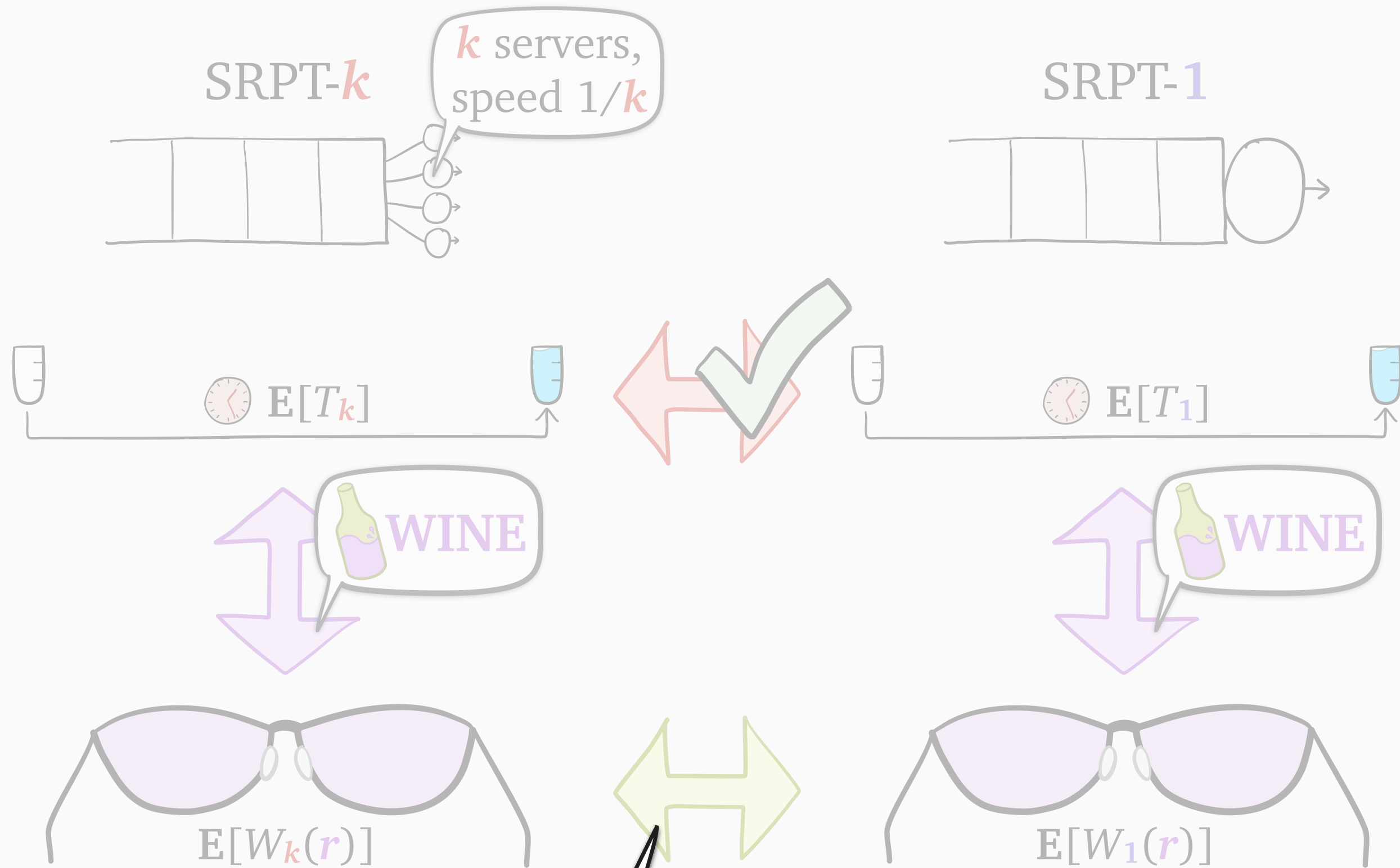


Lemma: r -work
decomposition



Lemma: r -work decomposition

$$E[W_k(r)] = E[W_1(r)] + \text{"}r\text{-work of } k - 1 \text{ jobs"}$$



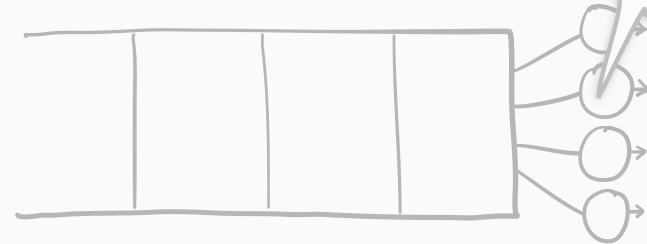
Lemma: r -work decomposition

$$E[W_k(r)] = E[W_1(r)] + \text{"}r\text{-work of } k - 1 \text{ jobs"}$$

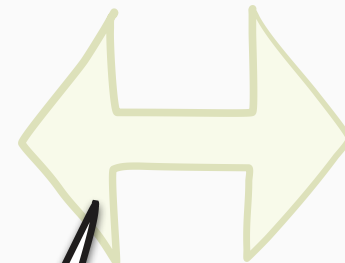
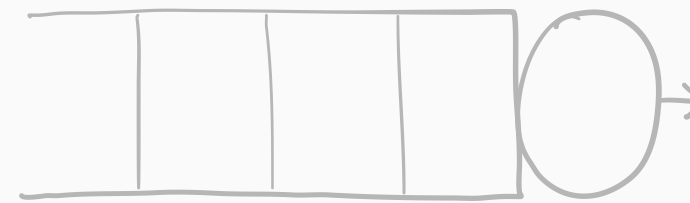
$$\leq E[W_1(r)] + (k - 1)r$$

SRPT- k

k servers,
speed $1/k$



SRPT-1



Lemma: r -work
decomposition

$$E[W_k(r)] = E[W_1(r)] + \text{"}r\text{-work of } k-1 \text{ jobs"}$$

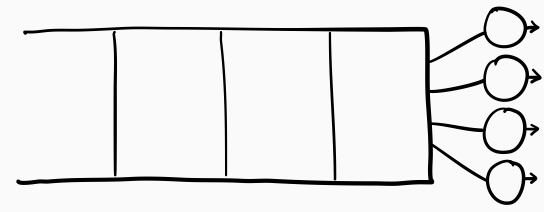
$$\leq E[W_1(r)] + (k-1)r$$

can improve



Part I

Handling job size uncertainty



Part II

Analyzing multiserver scheduling



Part III

Optimizing tail metrics

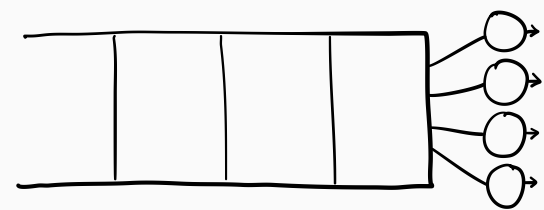


Part I

Handling job size uncertainty

Queueing for TCS

Use WINE to analyze SRPT-*k* with arbitrary release dates?



Part II

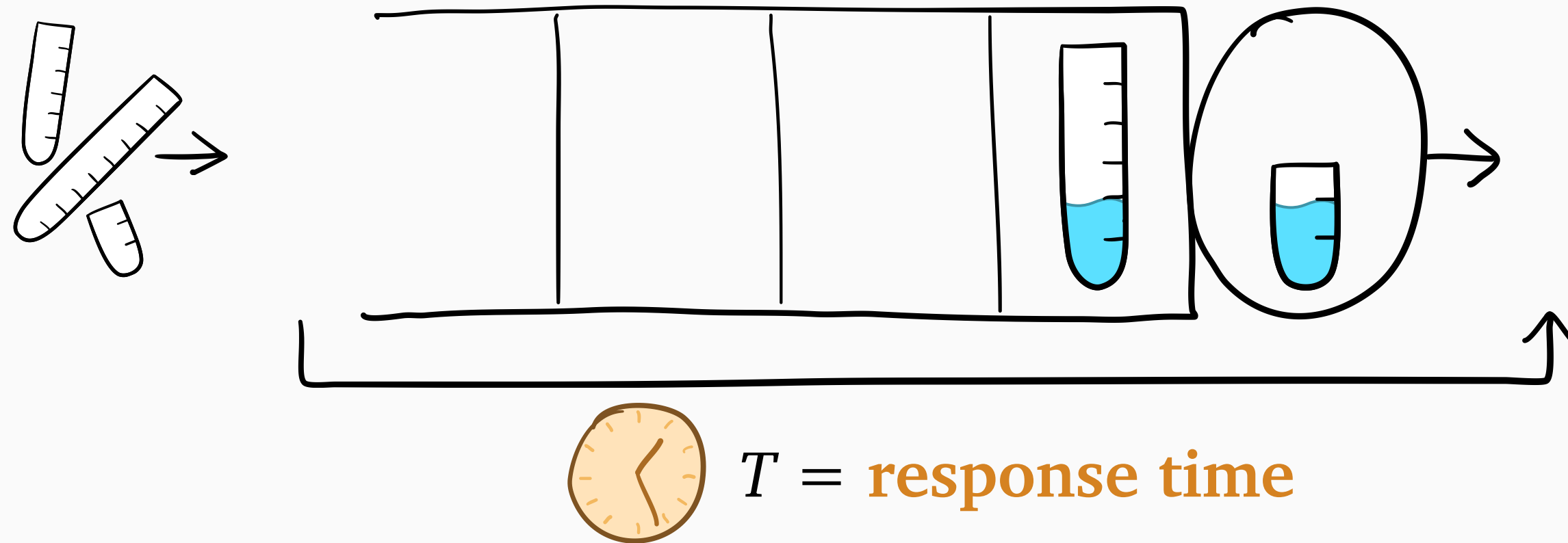
Analyzing multiserver scheduling



Part III

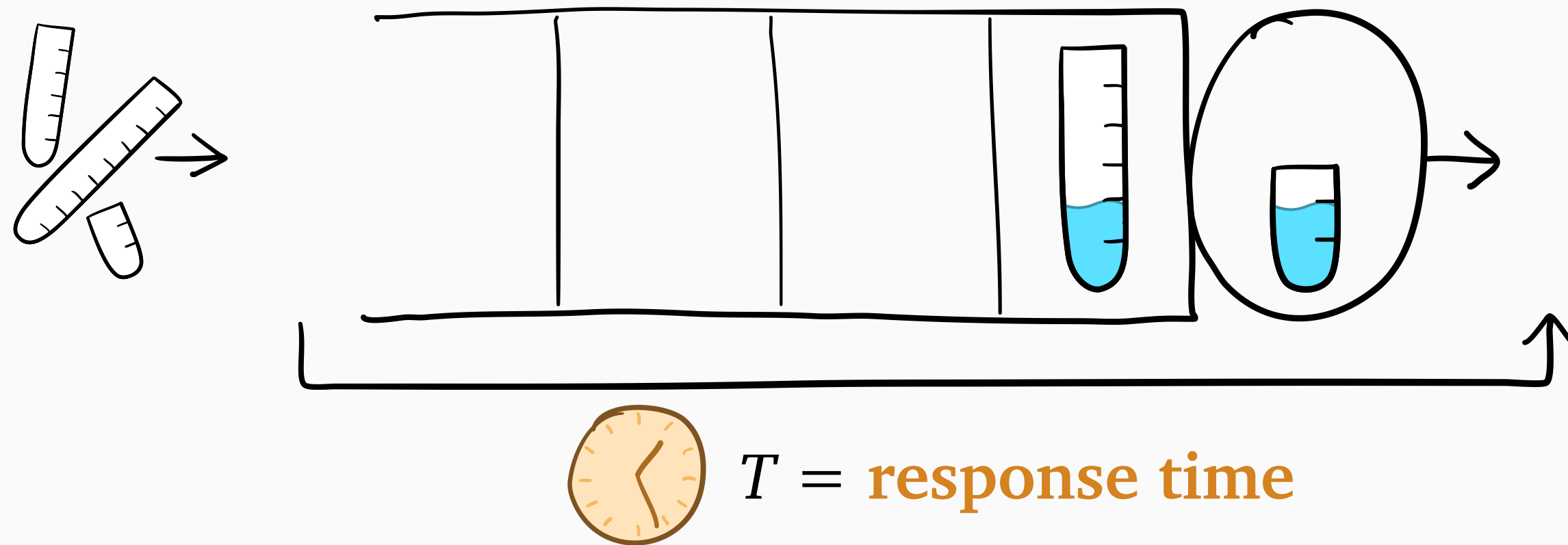
Optimizing tail metrics

Tail metrics



Tail metrics

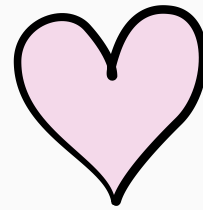
? Minimize $\begin{cases} \mathbf{P}[T > t]? \\ \mathbf{E}[(T - t)^+]? \\ \text{quantiles of } T? \end{cases}$



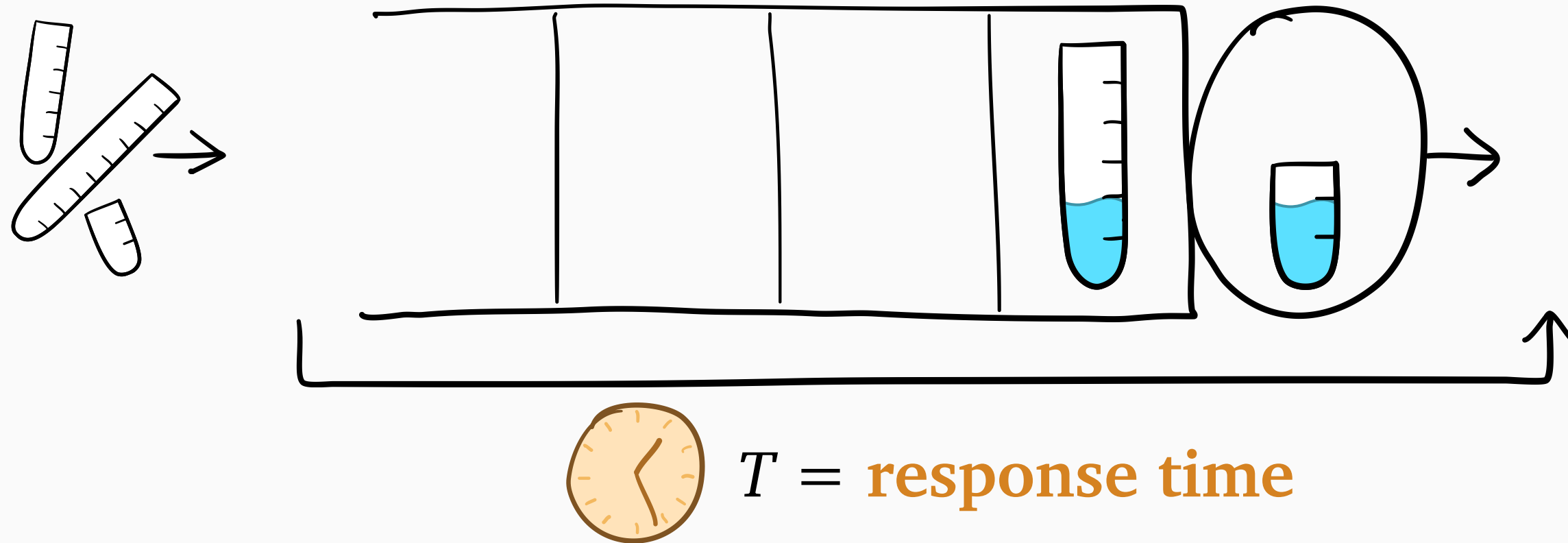
Tail metrics



Minimize $\begin{cases} \mathbf{P}[T > t]? \\ \mathbf{E}[(T - t)^+]? \\ \text{quantiles of } T? \end{cases}$



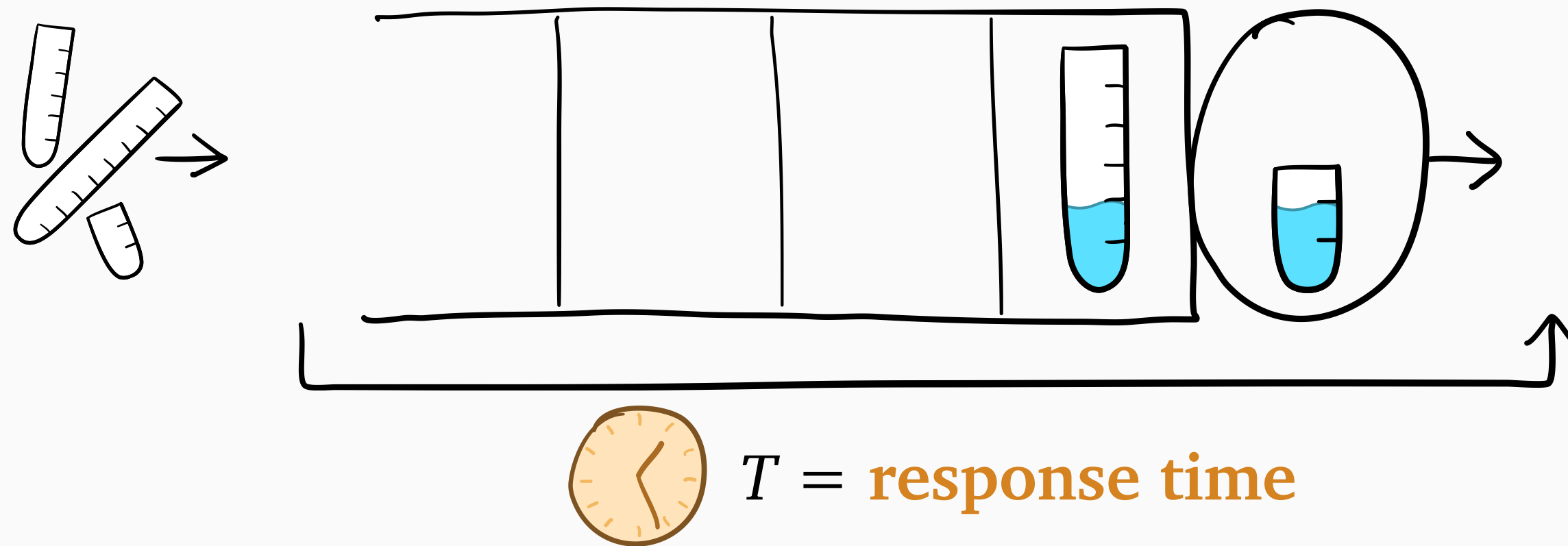
Practice: important



Tail metrics

? Minimize $\begin{cases} \mathbf{P}[T > t]? \\ \mathbf{E}[(T - t)^+]? \\ \text{quantiles of } T? \end{cases}$

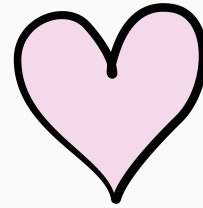
♥ Practice: important
⚠ Theory: very hard



Tail metrics



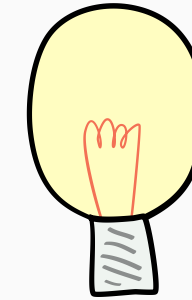
Minimize $\begin{cases} \mathbf{P}[T > t]? \\ \mathbf{E}[(T - t)^+]? \\ \text{quantiles of } T? \end{cases}$



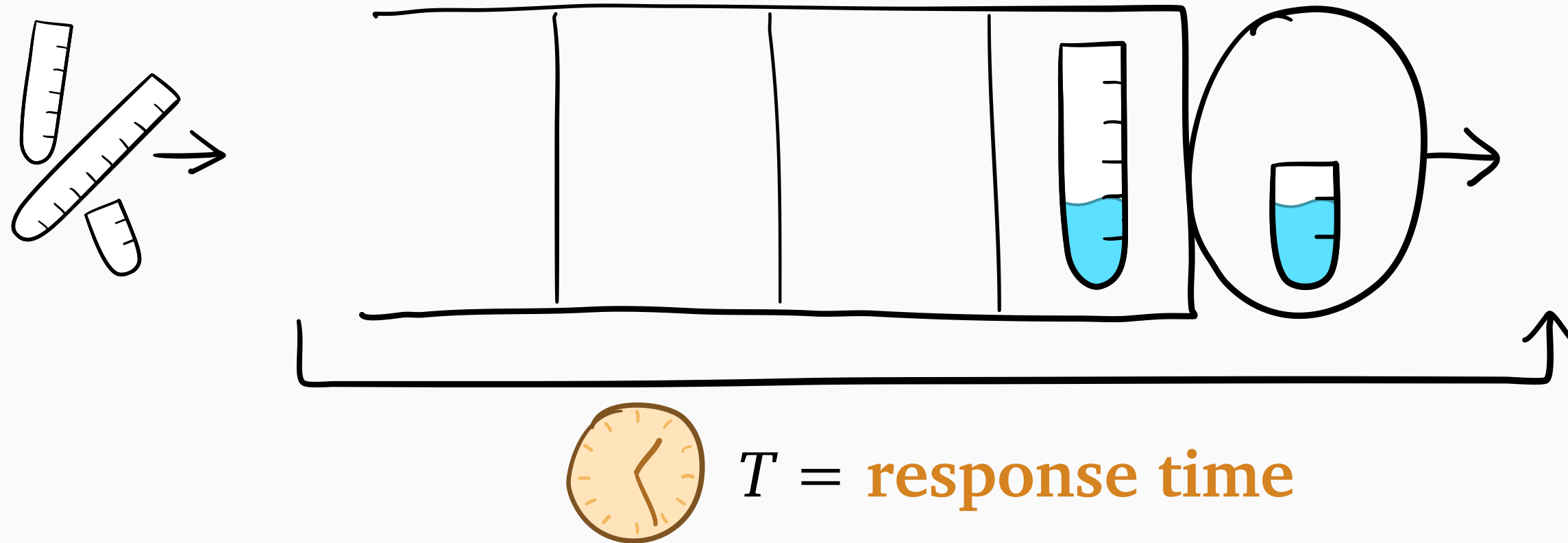
Practice: important



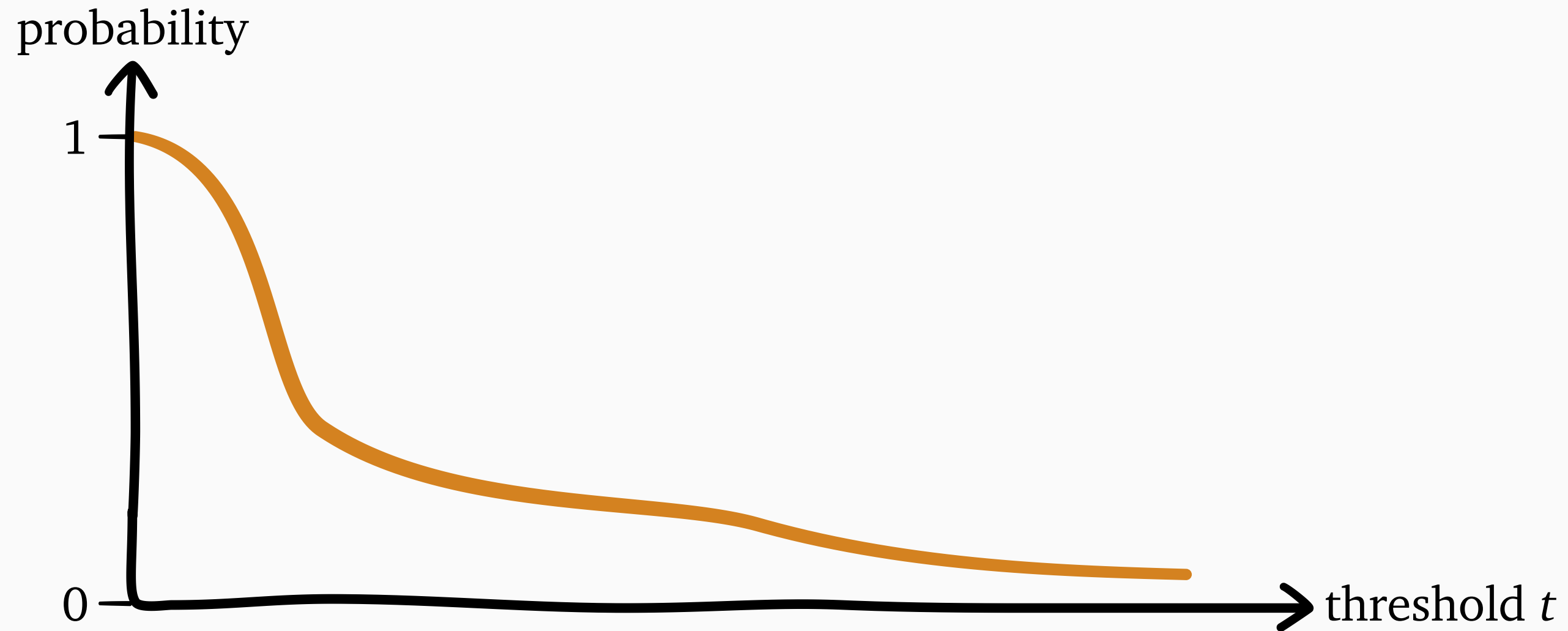
Theory: very hard



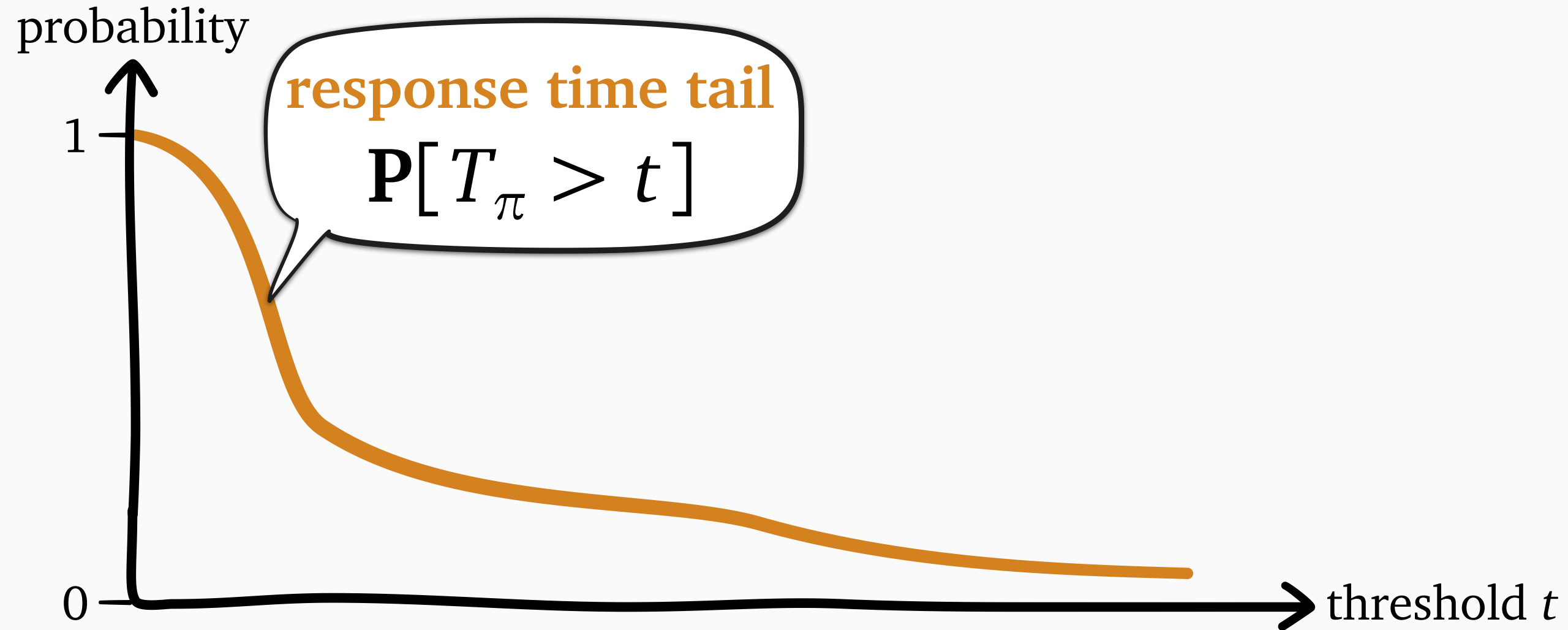
Tractable:
study $t \rightarrow \infty$
asymptotics



Asymptotic response time tail



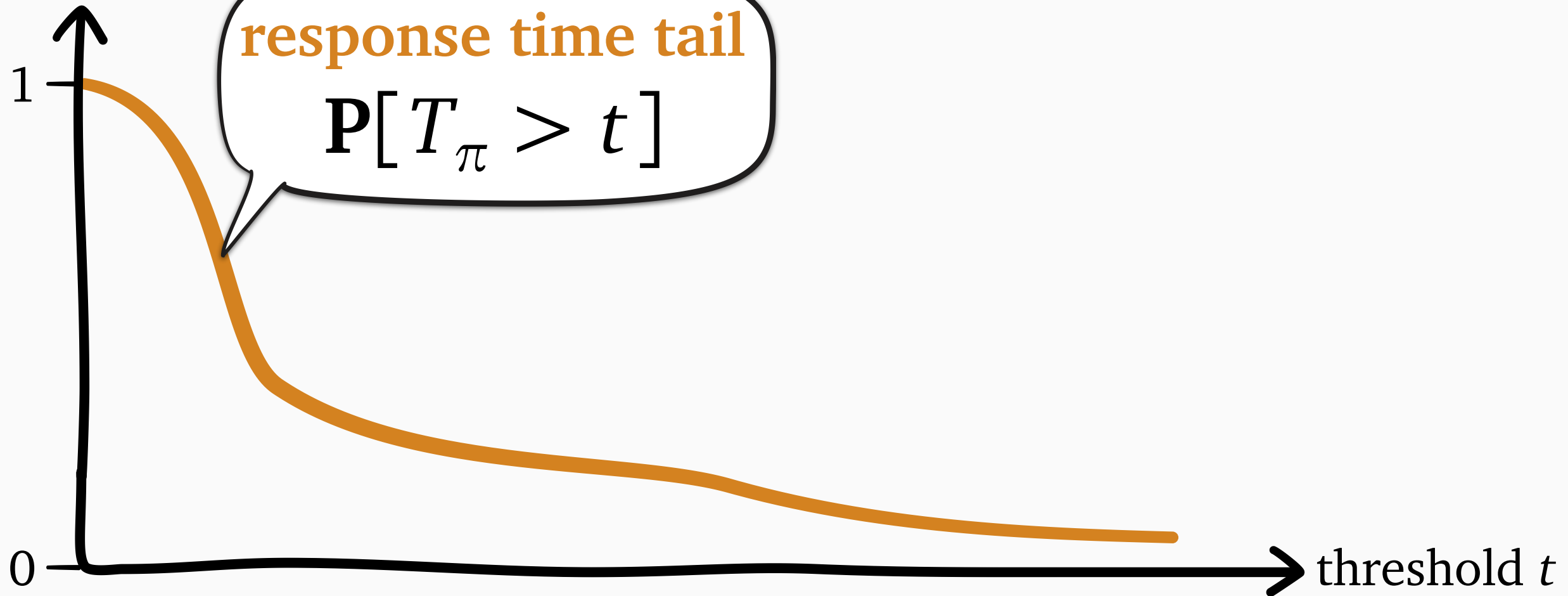
Asymptotic response time tail



Asymptotic response time tail

depends on
policy π

probability



Asymptotic response time tail

depends on
policy π

probability

response time tail

$$\mathbf{P}[T_{\pi} > t]$$

1

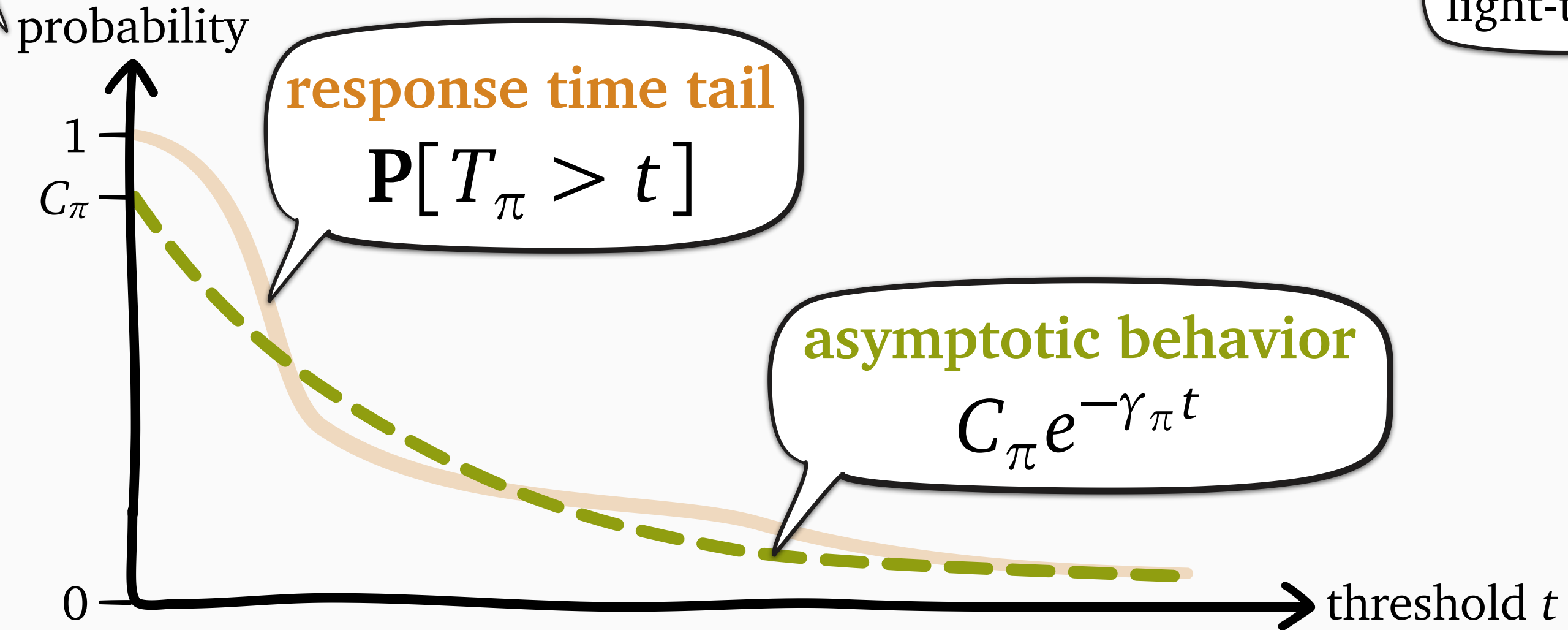
0

threshold t

Asymptotic response time tail

depends on
policy π

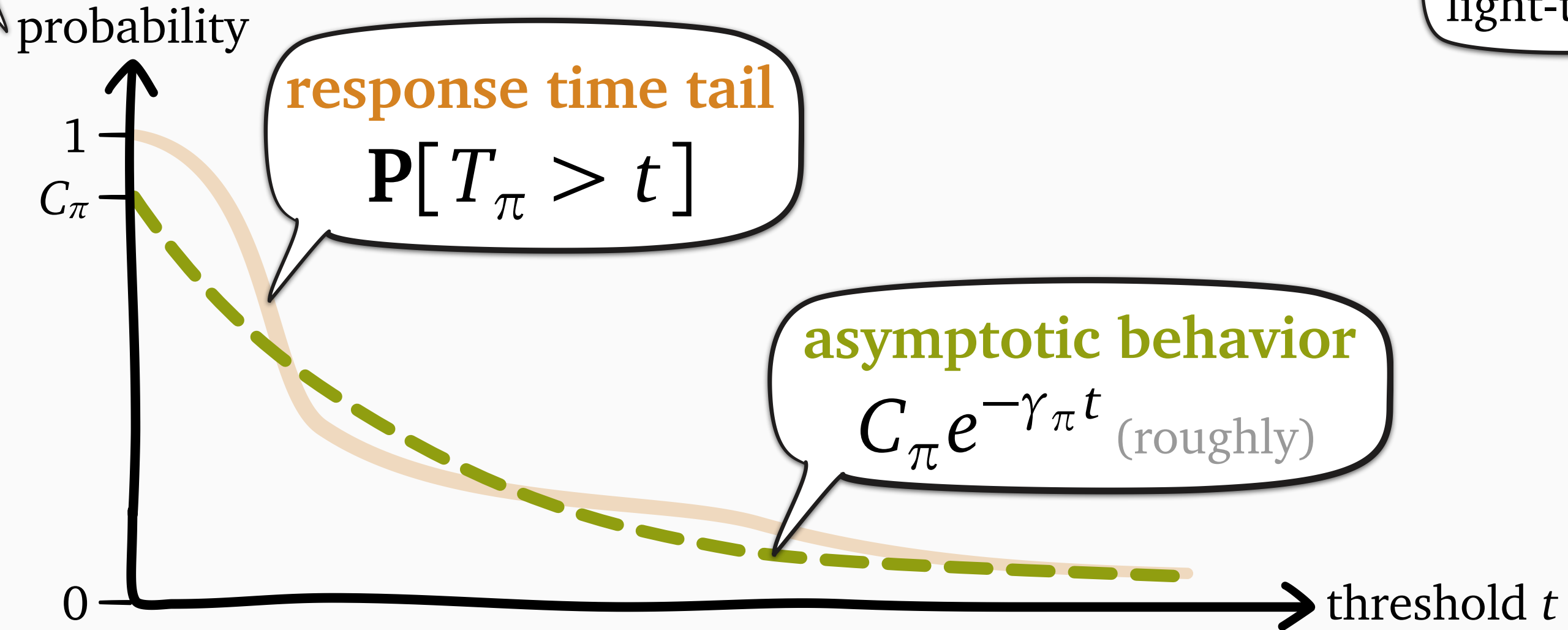
when S is
light-tailed



Asymptotic response time tail

depends on
policy π

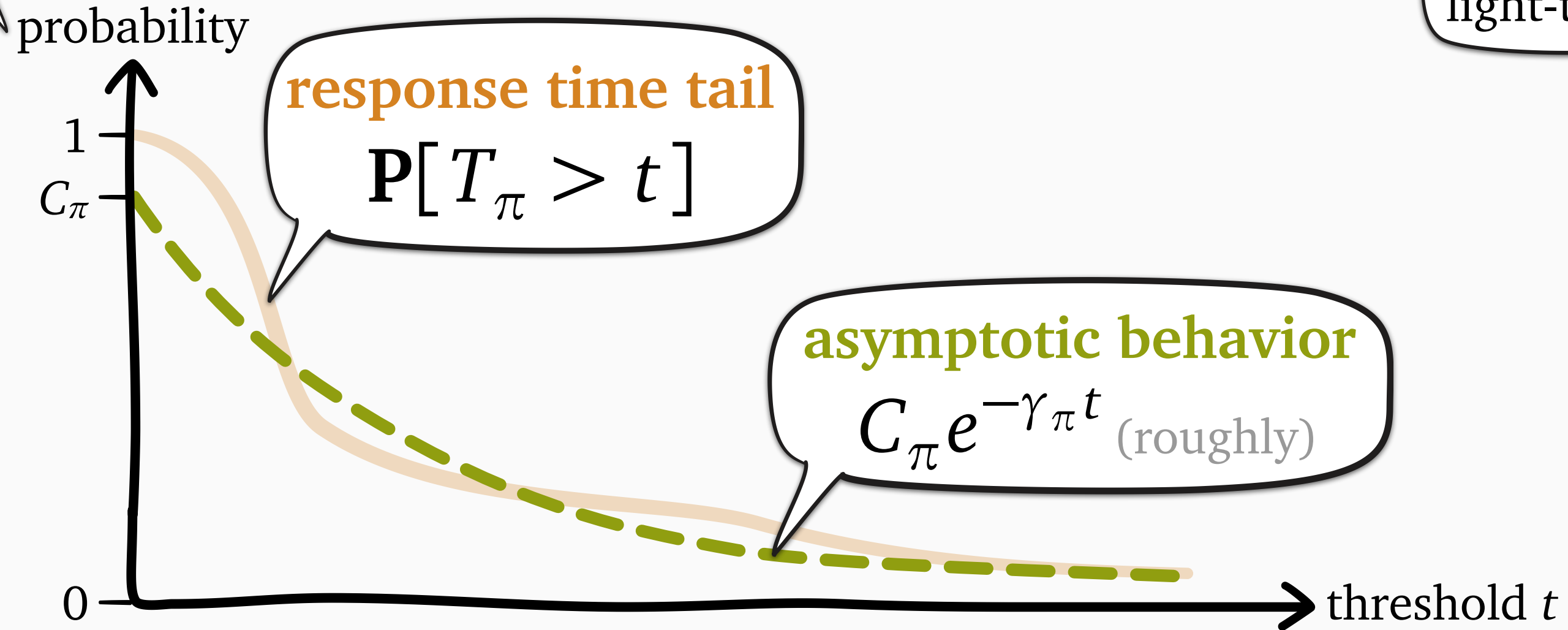
when S is
light-tailed



Asymptotic response time tail

depends on
policy π

when S is
light-tailed



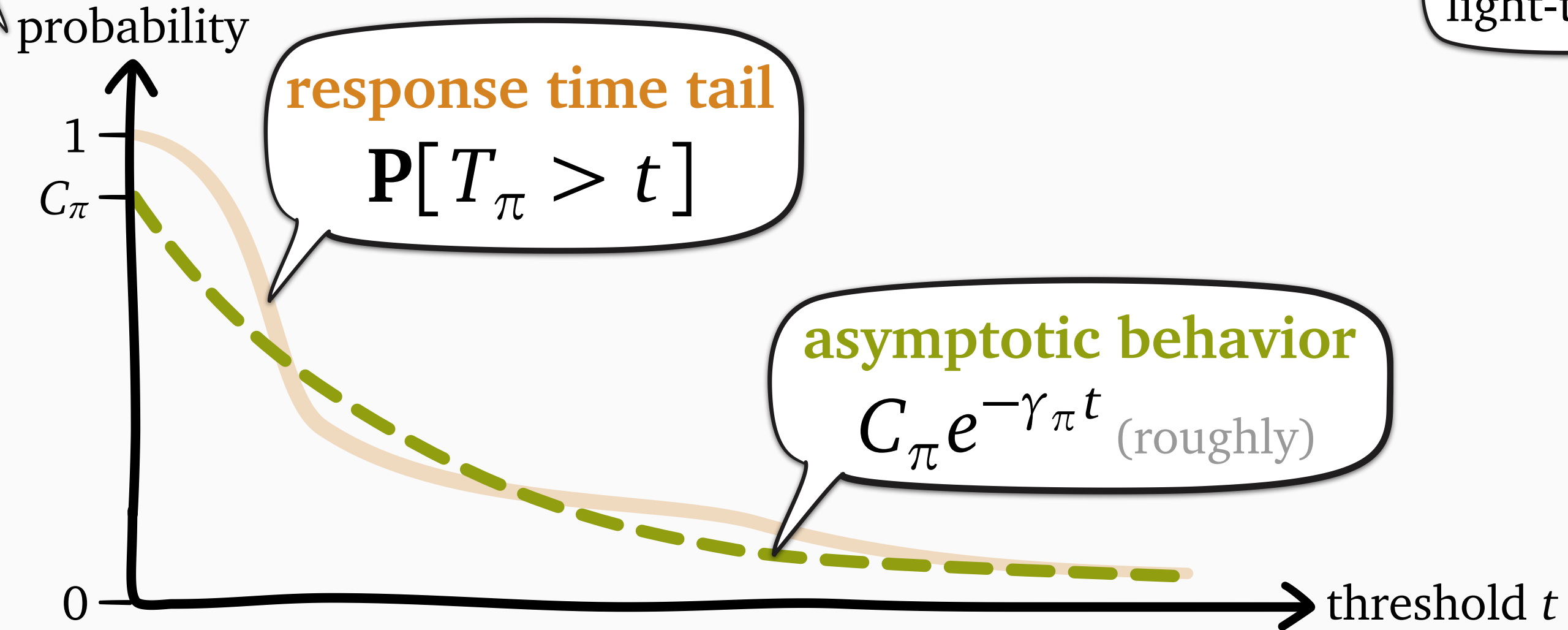
$\gamma_\pi = \text{decay rate of } \pi$

$C_\pi = \text{tail constant of } \pi$

Asymptotic response time tail

depends on
policy π

when S is
light-tailed



Weak optimality: ←

optimal γ_π

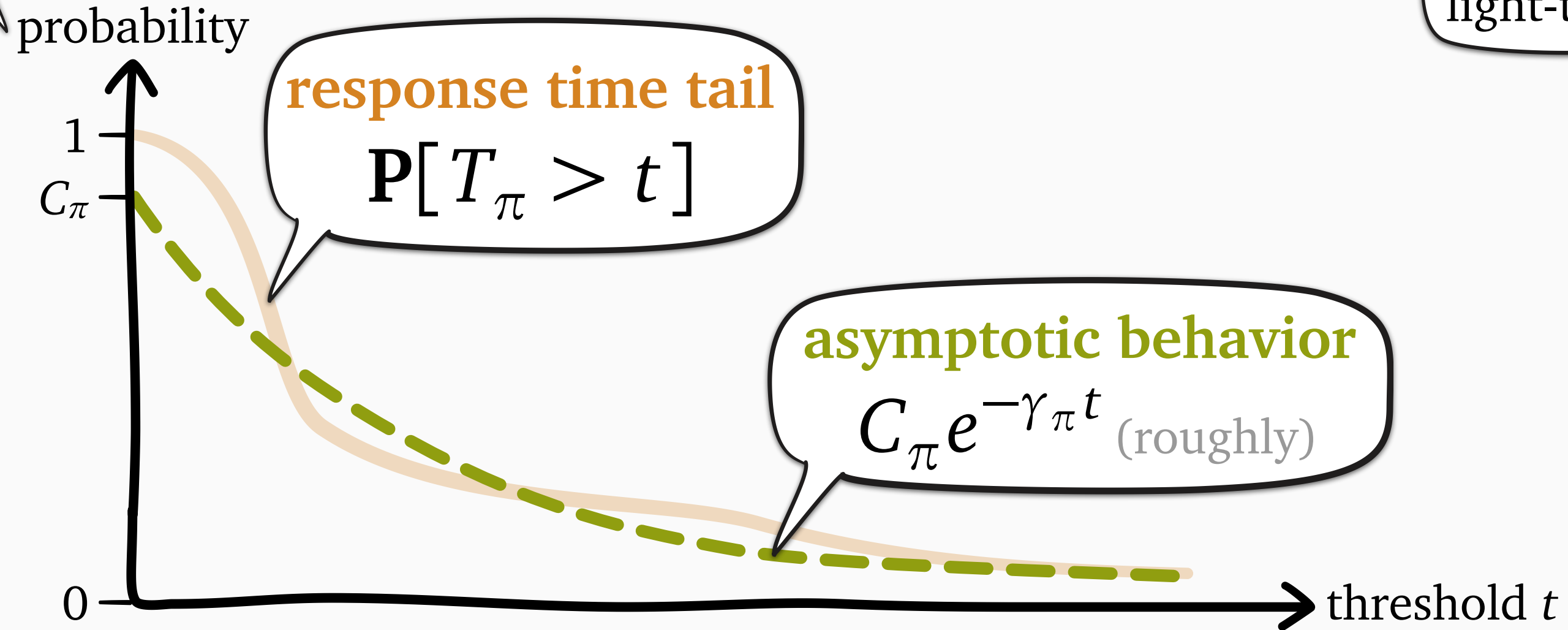
$\gamma_\pi = \text{decay rate of } \pi$

$C_\pi = \text{tail constant of } \pi$

Asymptotic response time tail

depends on
policy π

when S is
light-tailed



Weak optimality: ←
optimal γ_π

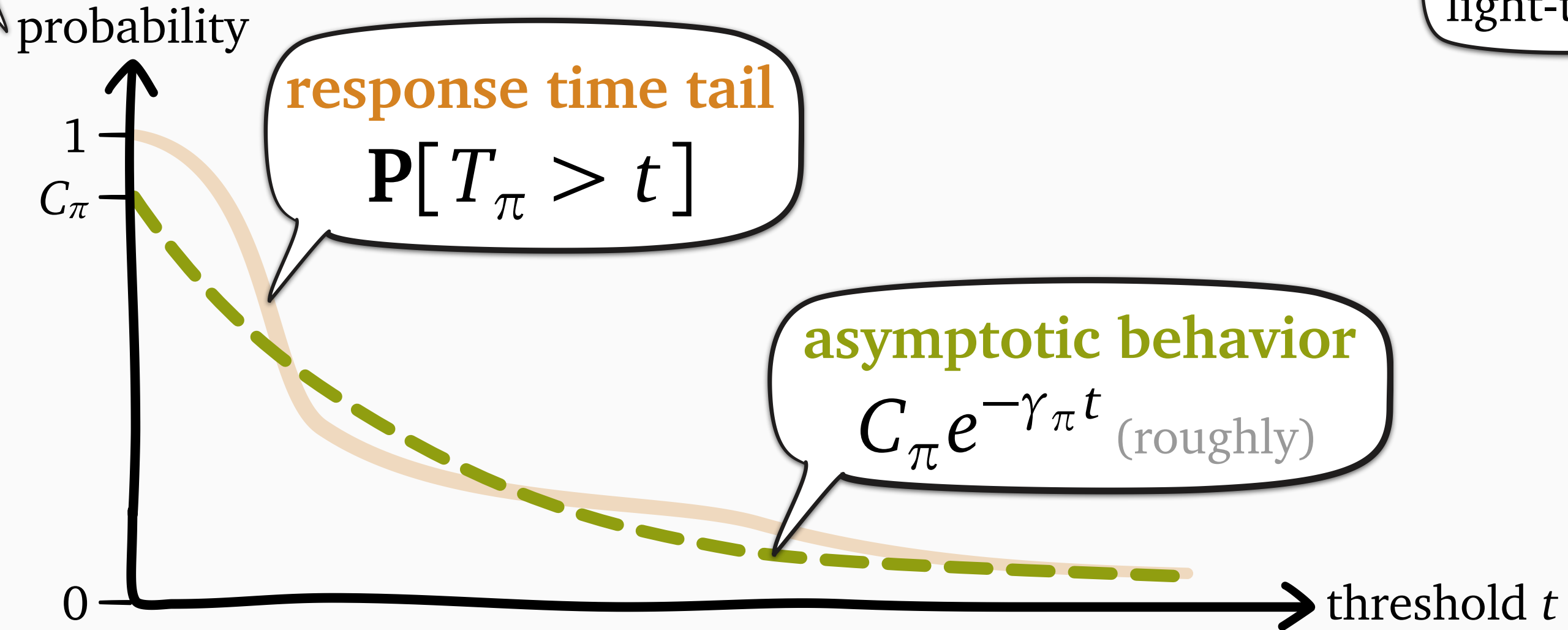
$\gamma_\pi = \text{decay rate of } \pi$
 $C_\pi = \text{tail constant of } \pi$

→ **Strong optimality:**
optimal γ_π and C_π

Asymptotic response time tail

depends on
policy π

when S is
light-tailed



Weak optimality: ←
optimal γ_π

$\gamma_\pi = \text{decay rate of } \pi$
 $C_\pi = \text{tail constant of } \pi$

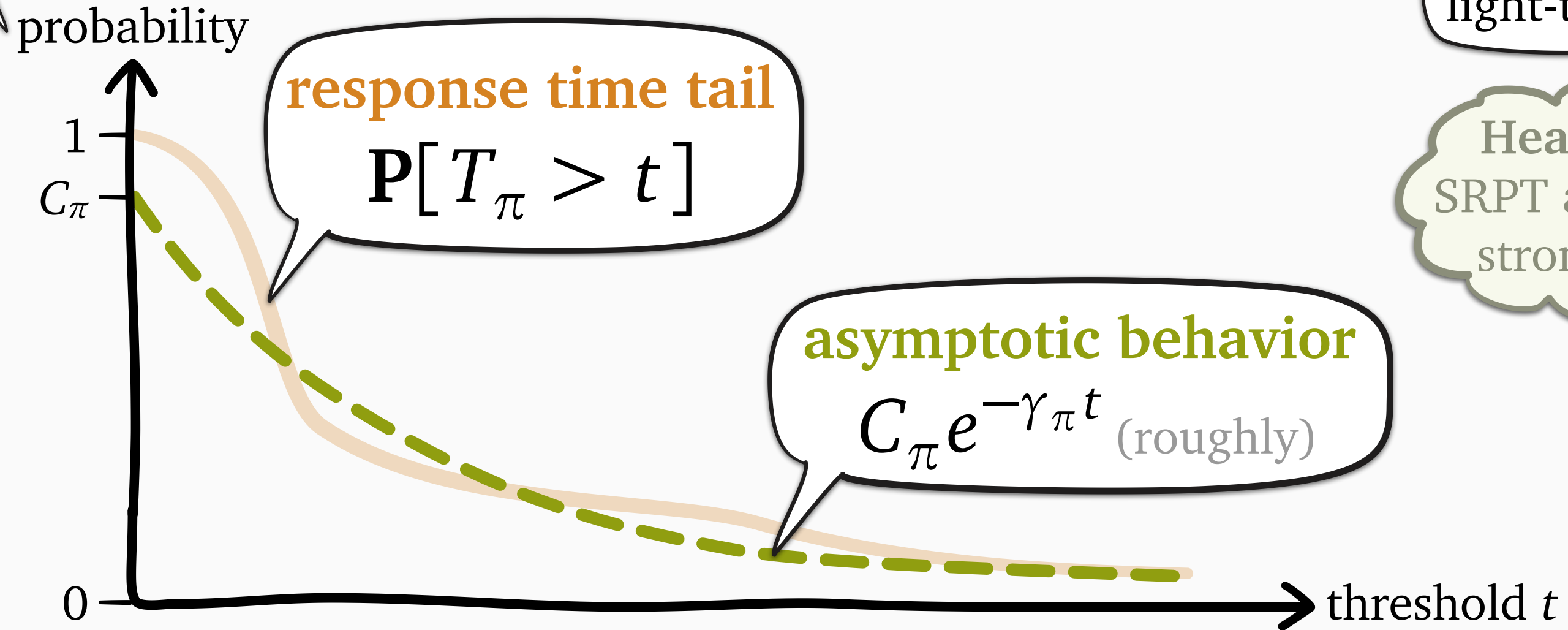
→ **Strong optimality:**
optimal γ_π and C_π
(roughly)

Asymptotic response time tail

depends on
policy π

when S is
light-tailed

Heavy-tailed S :
SRPT and others are
strongly optimal

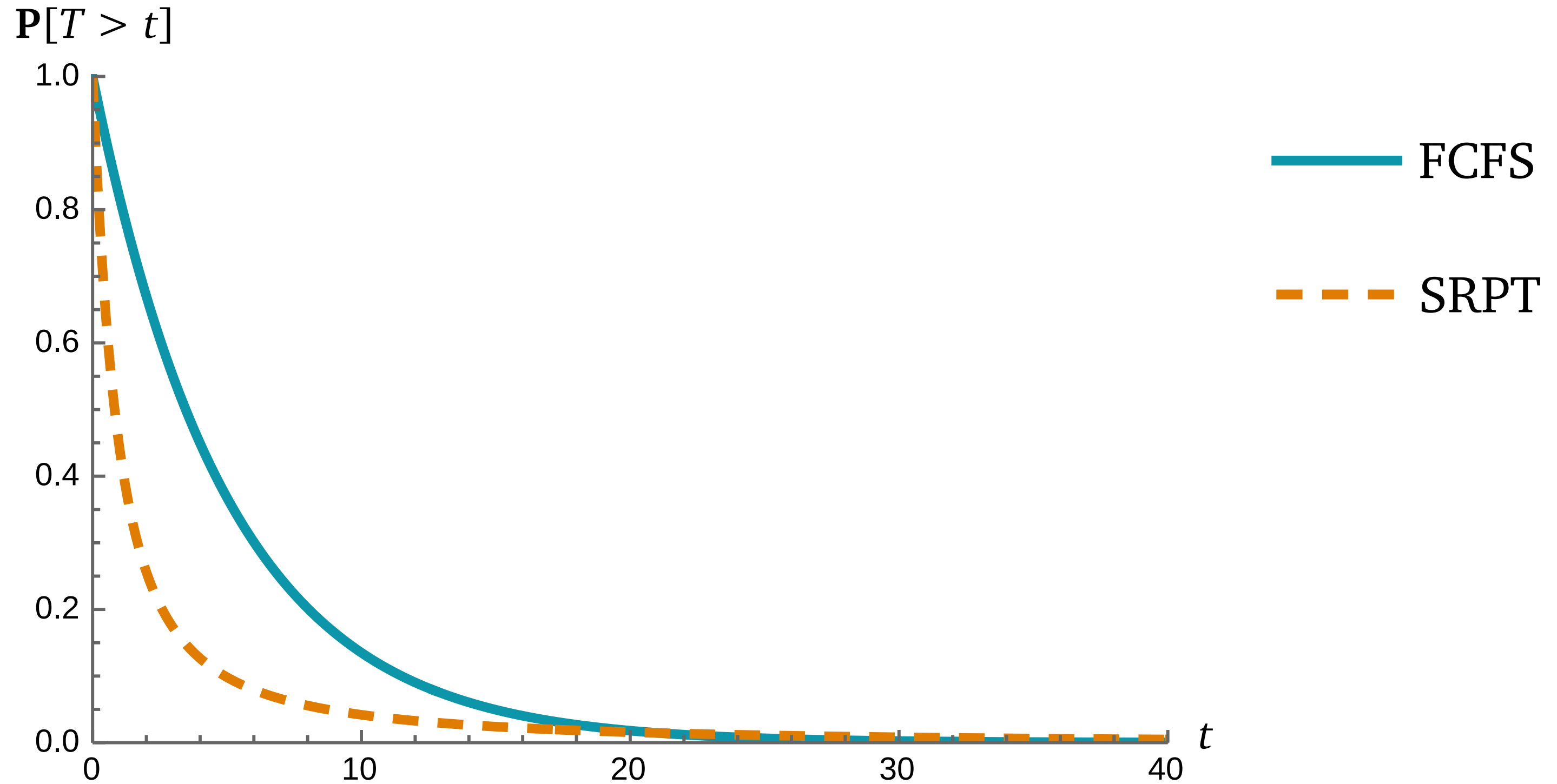


Weak optimality: ←
optimal γ_π

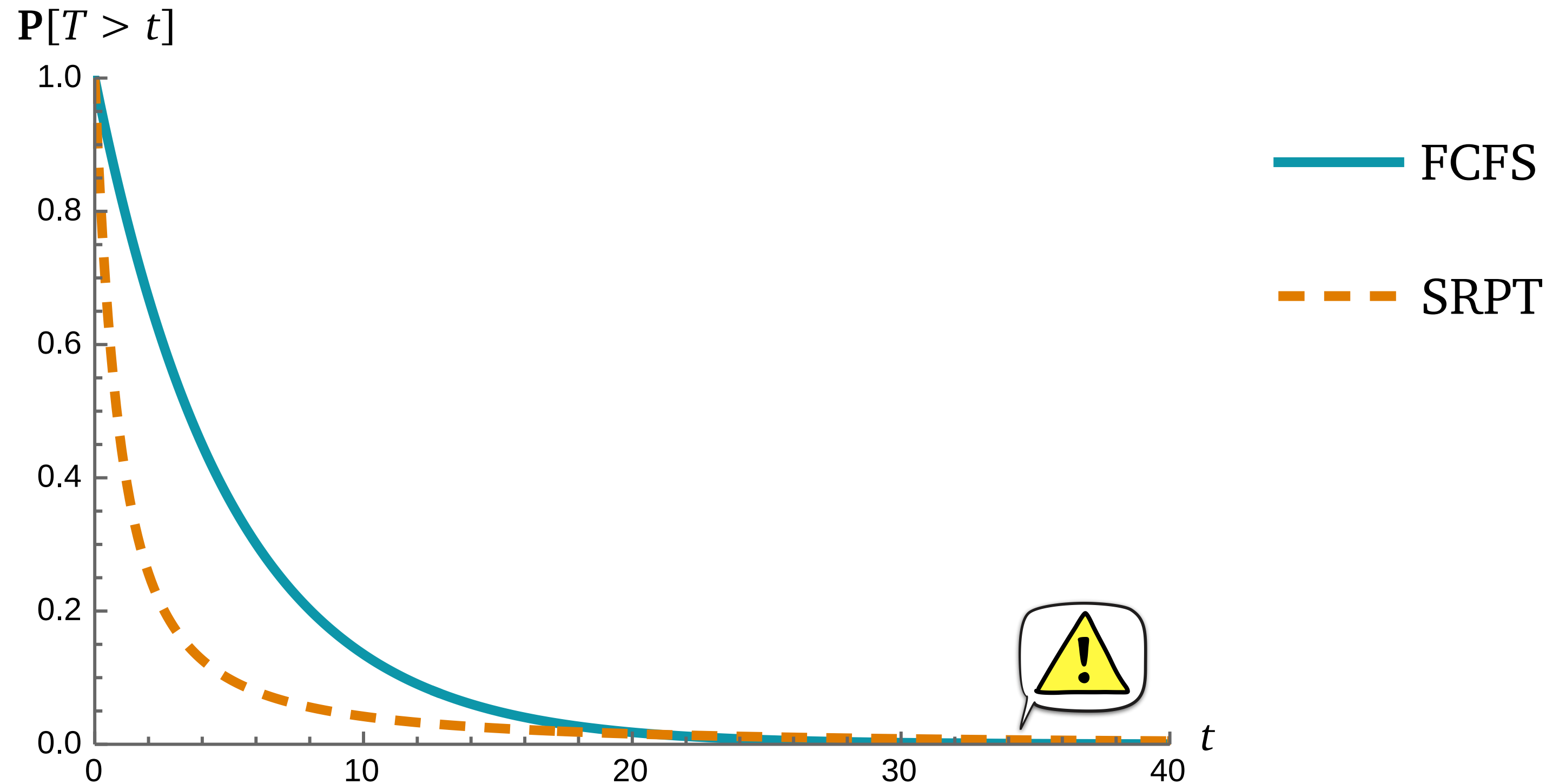
$\gamma_\pi = \text{decay rate of } \pi$
 $C_\pi = \text{tail constant of } \pi$

→ **Strong optimality:**
optimal γ_π and C_π
(roughly)

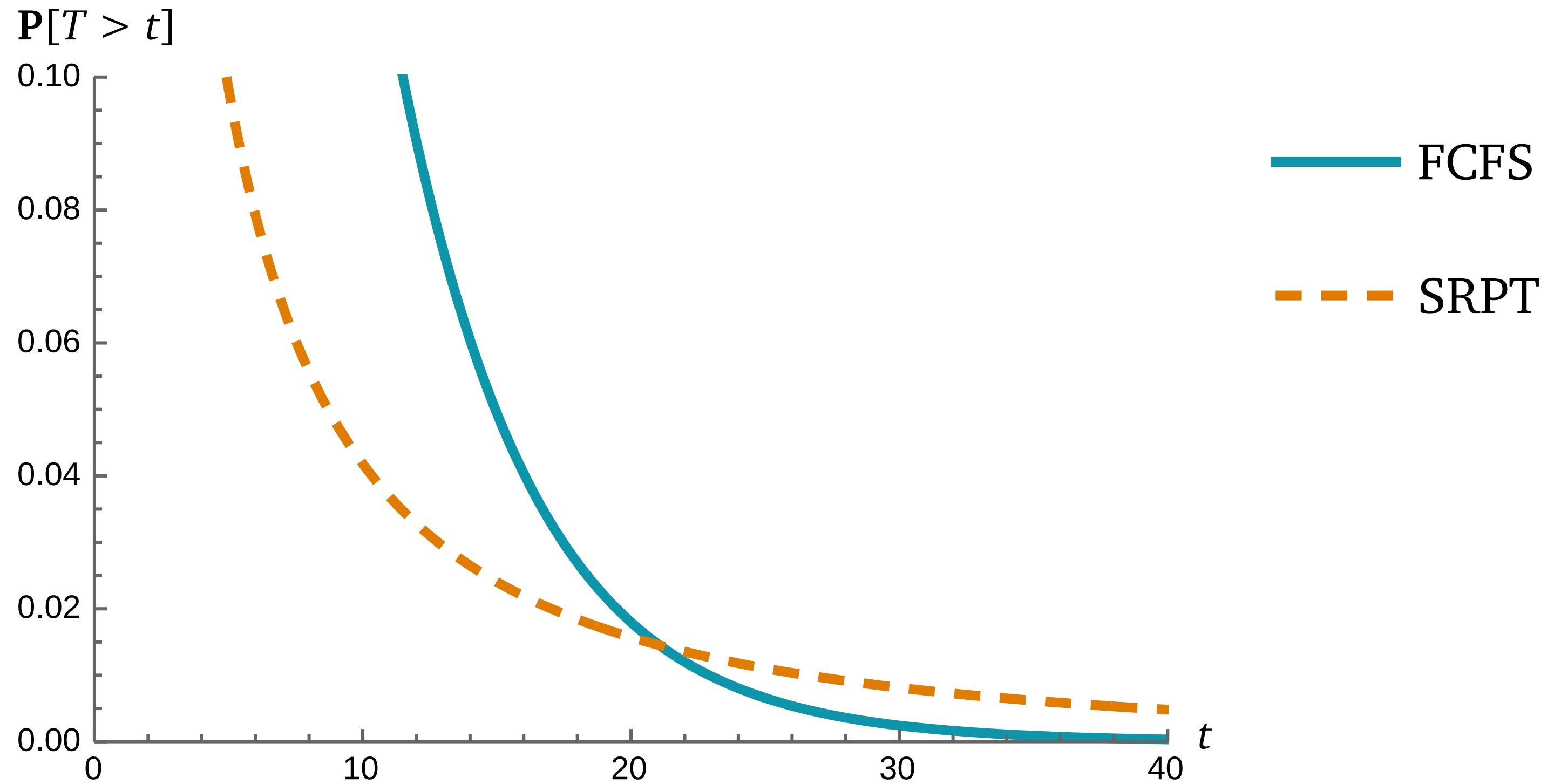
Optimizing the decay rate γ



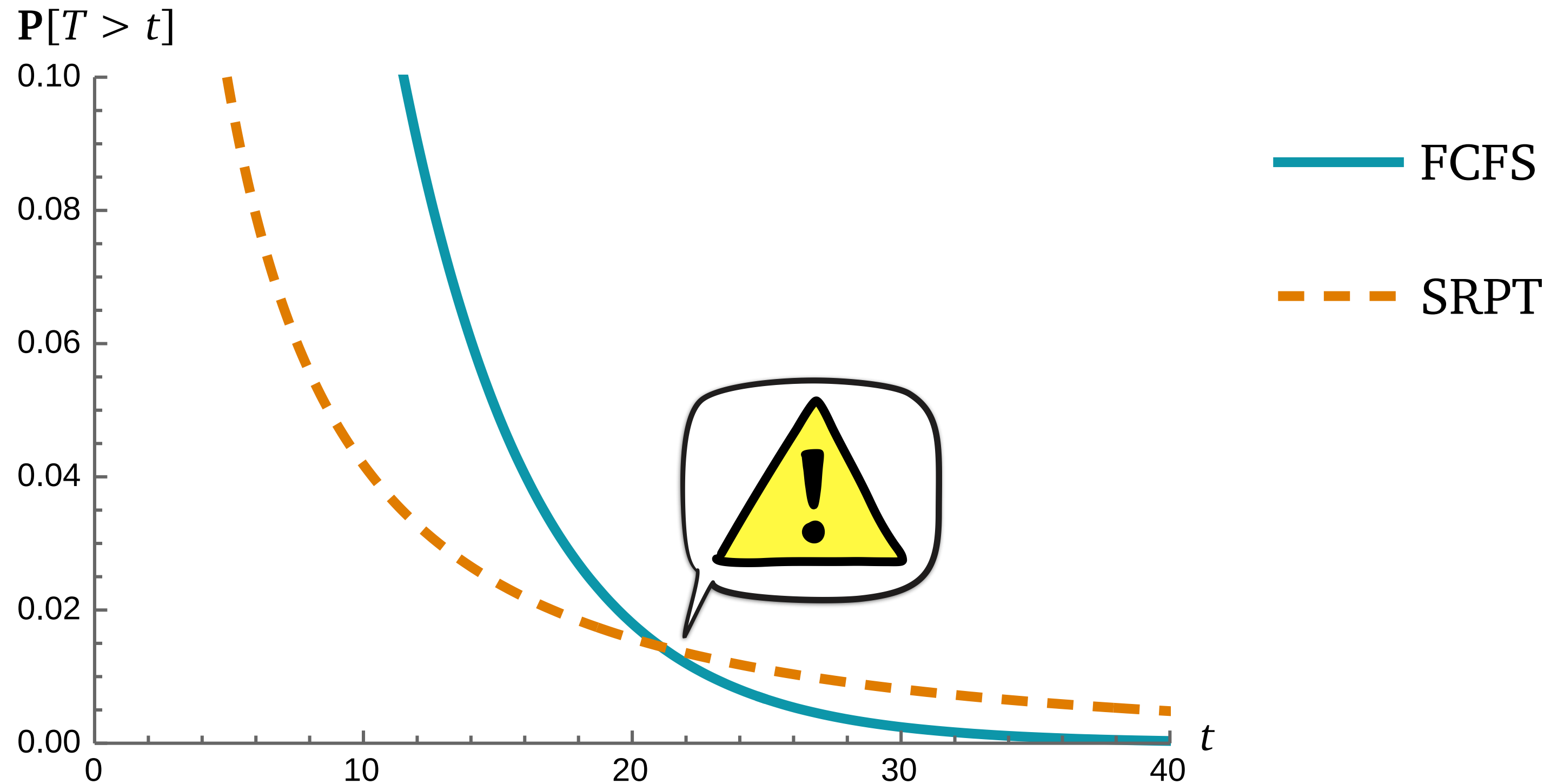
Optimizing the decay rate γ



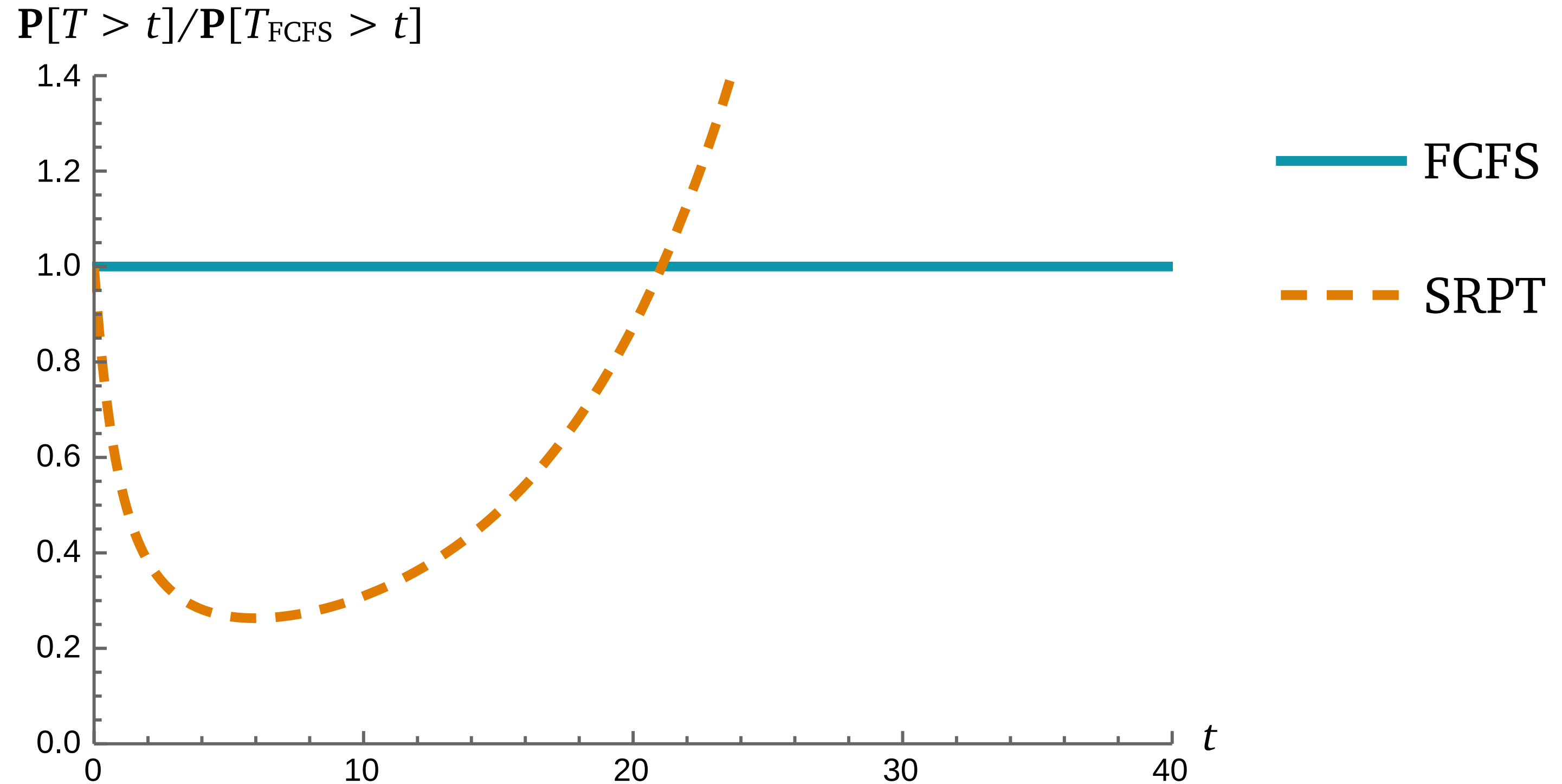
Optimizing the decay rate γ



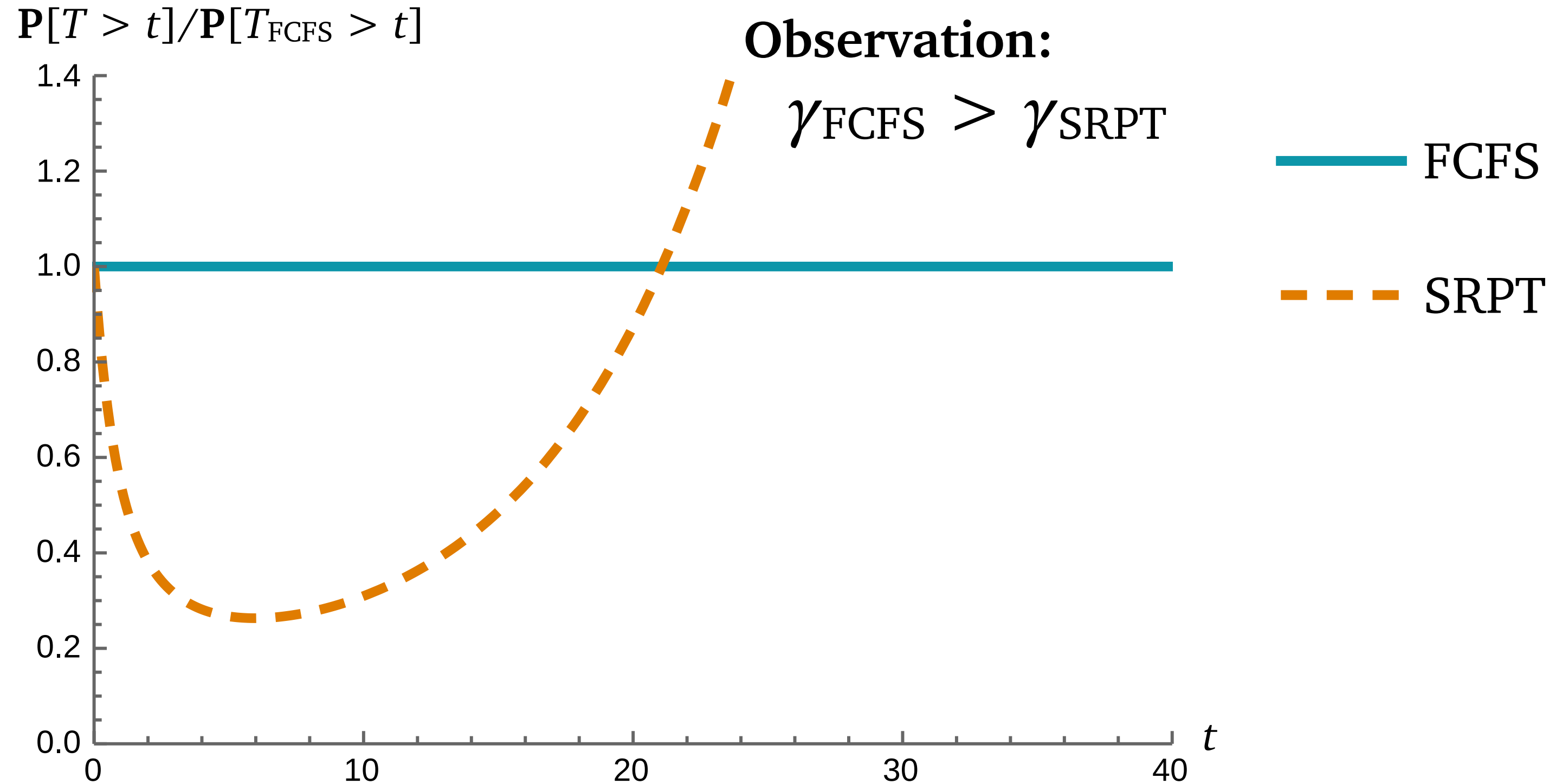
Optimizing the decay rate γ



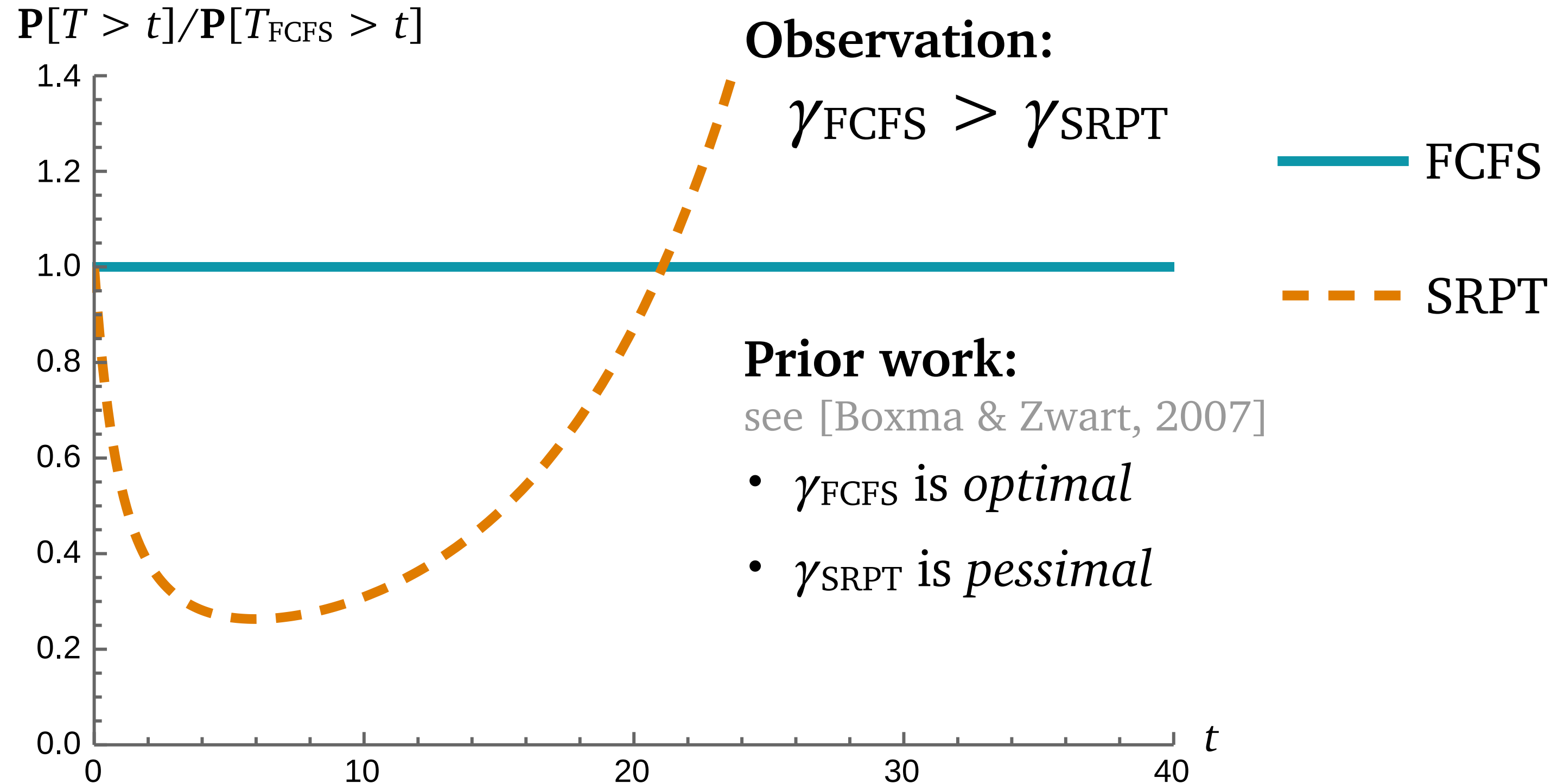
Optimizing the decay rate γ



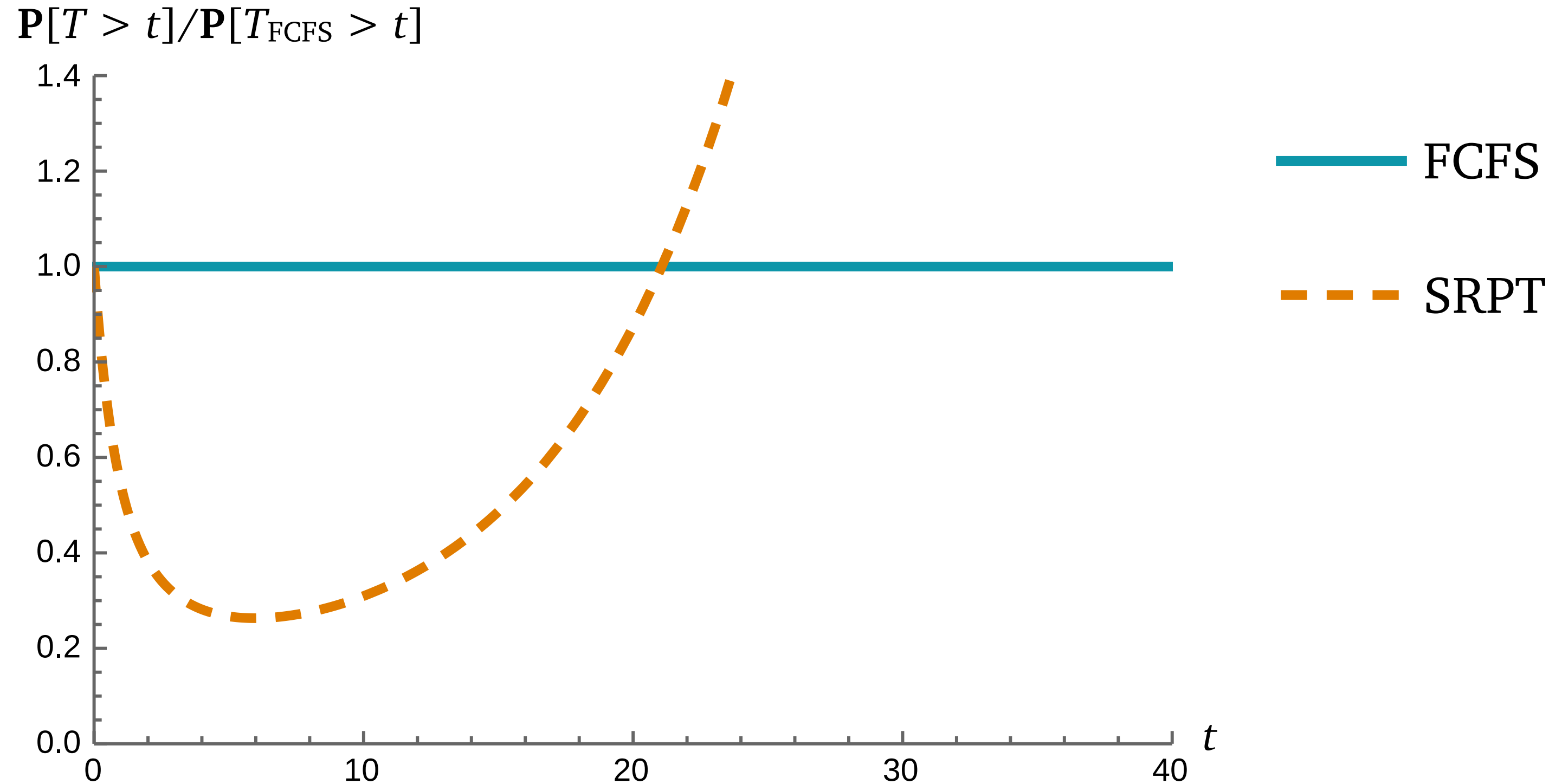
Optimizing the decay rate γ



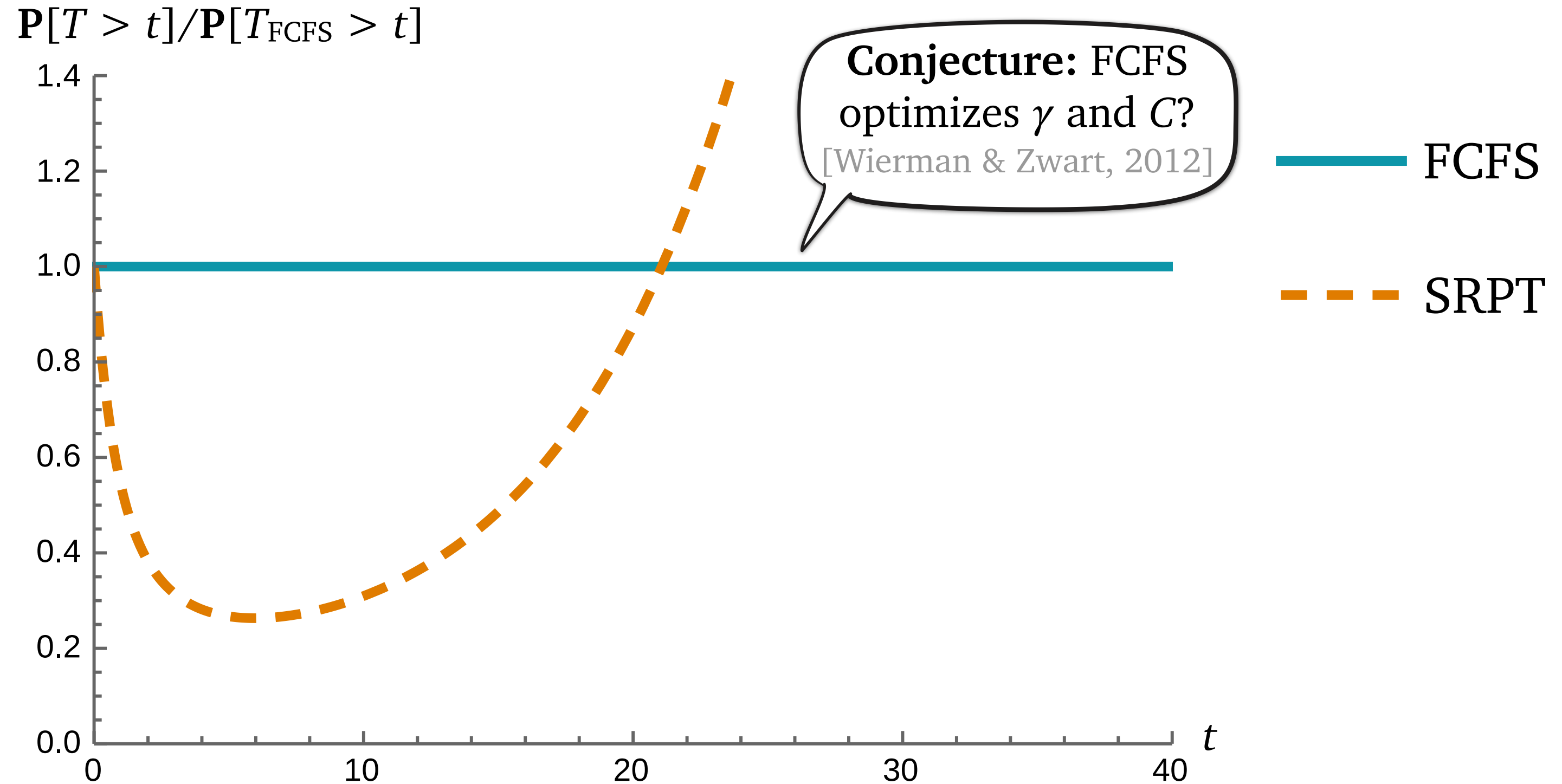
Optimizing the decay rate γ



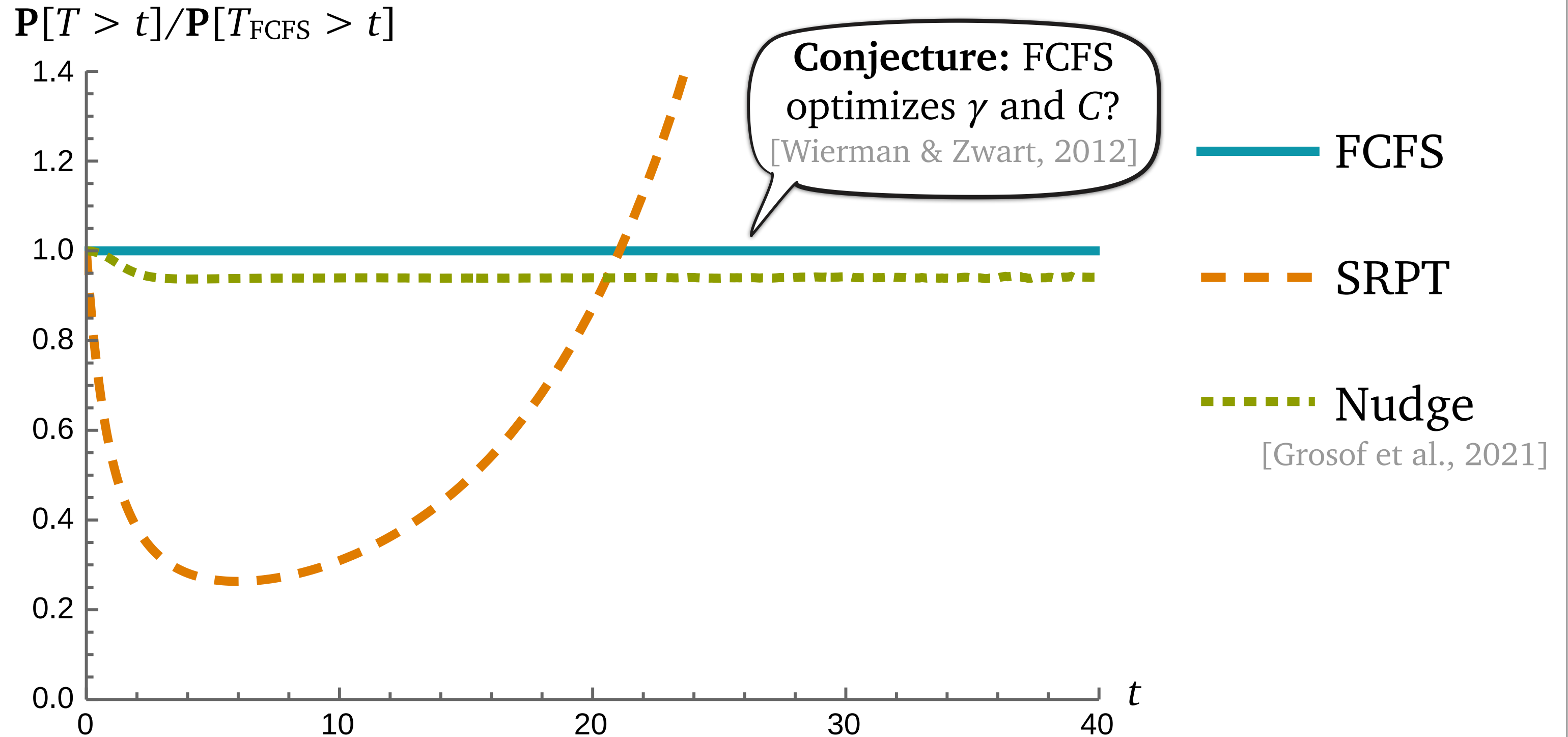
Optimizing the tail constant C



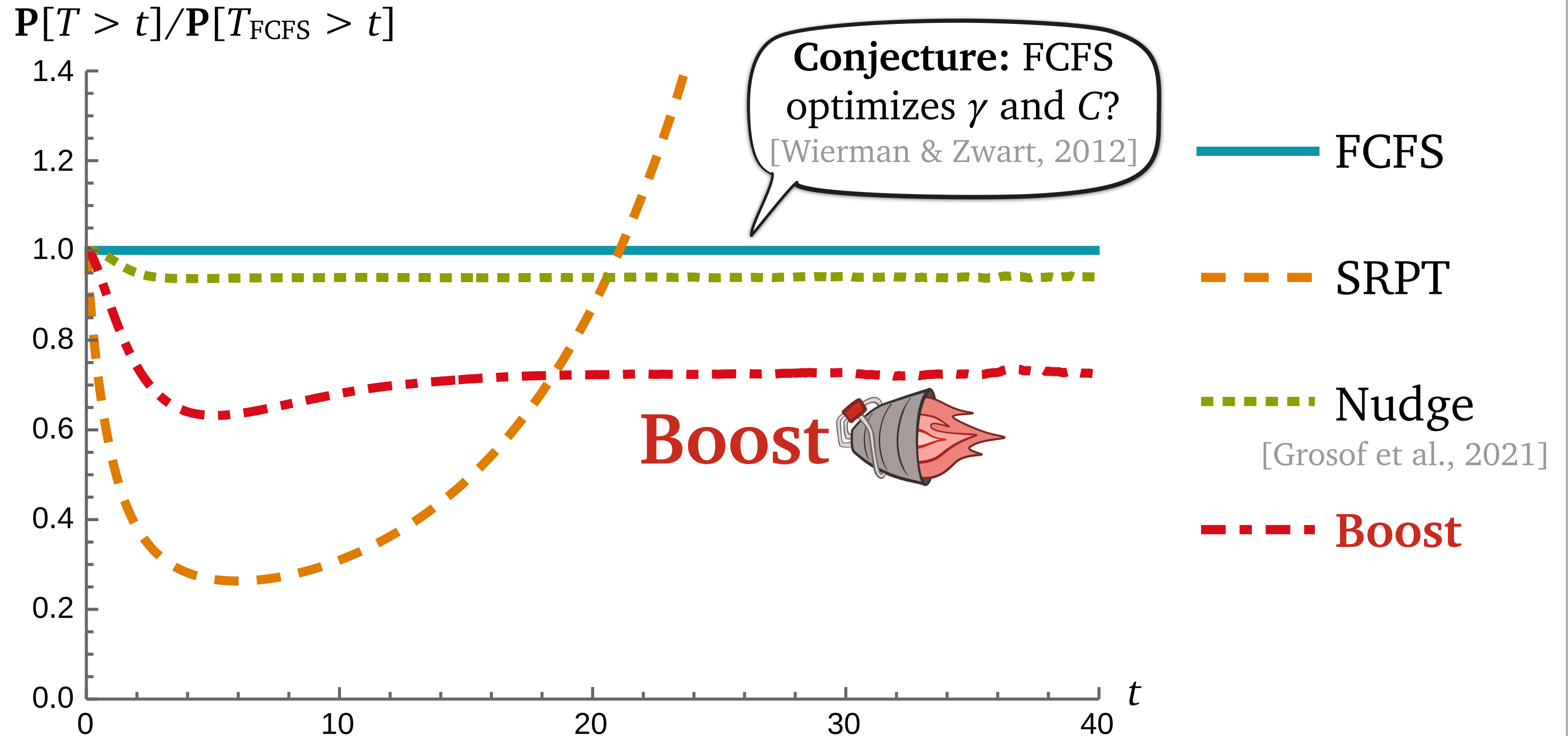
Optimizing the tail constant C



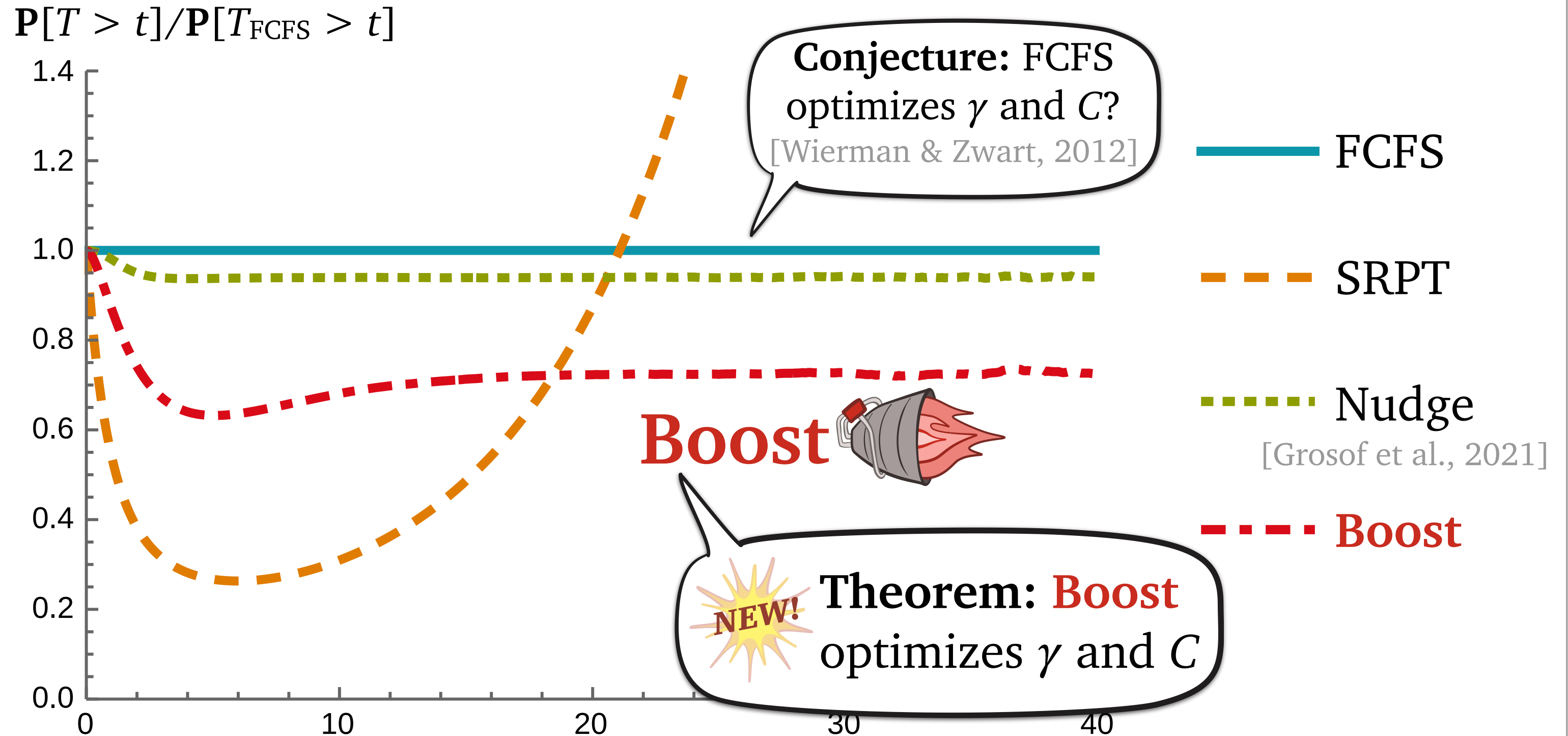
Optimizing the tail constant C



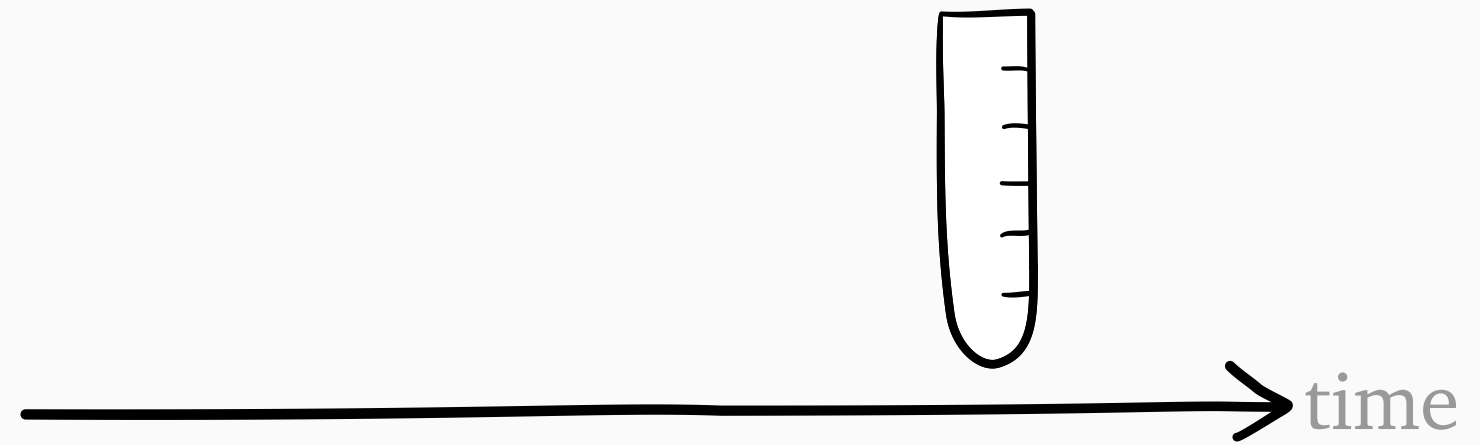
Optimizing the tail constant C



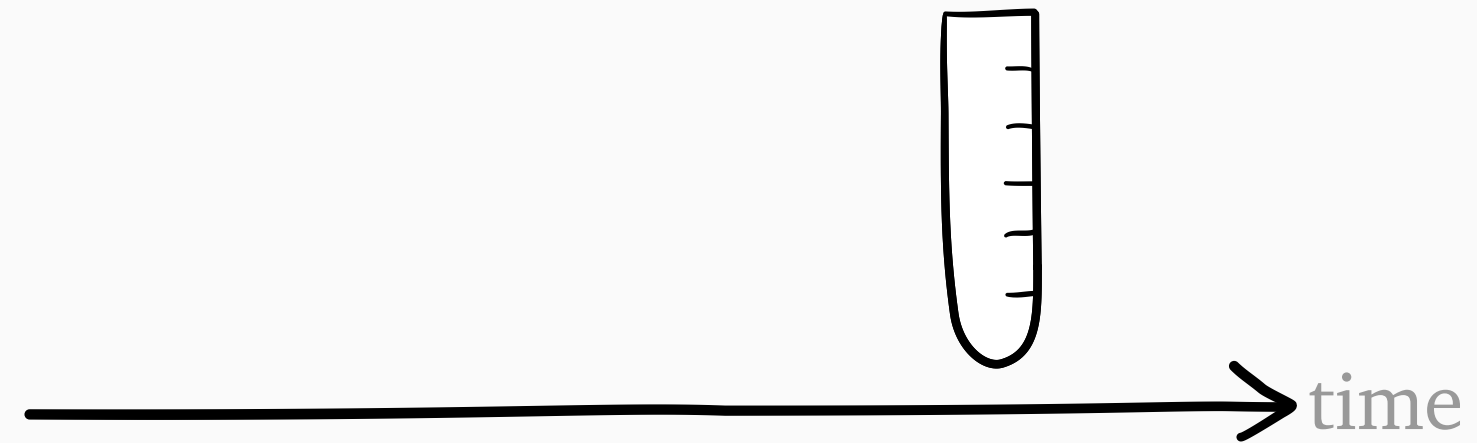
Optimizing the tail constant C



How **Boost** works

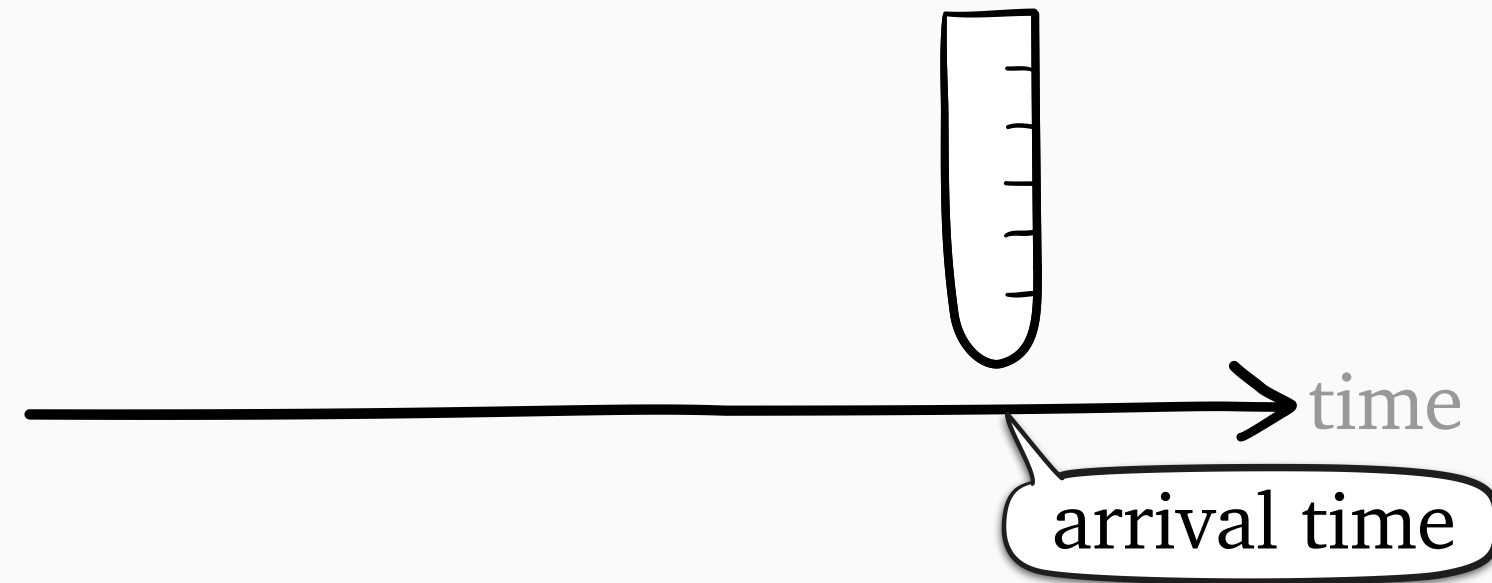


How **Boost** works



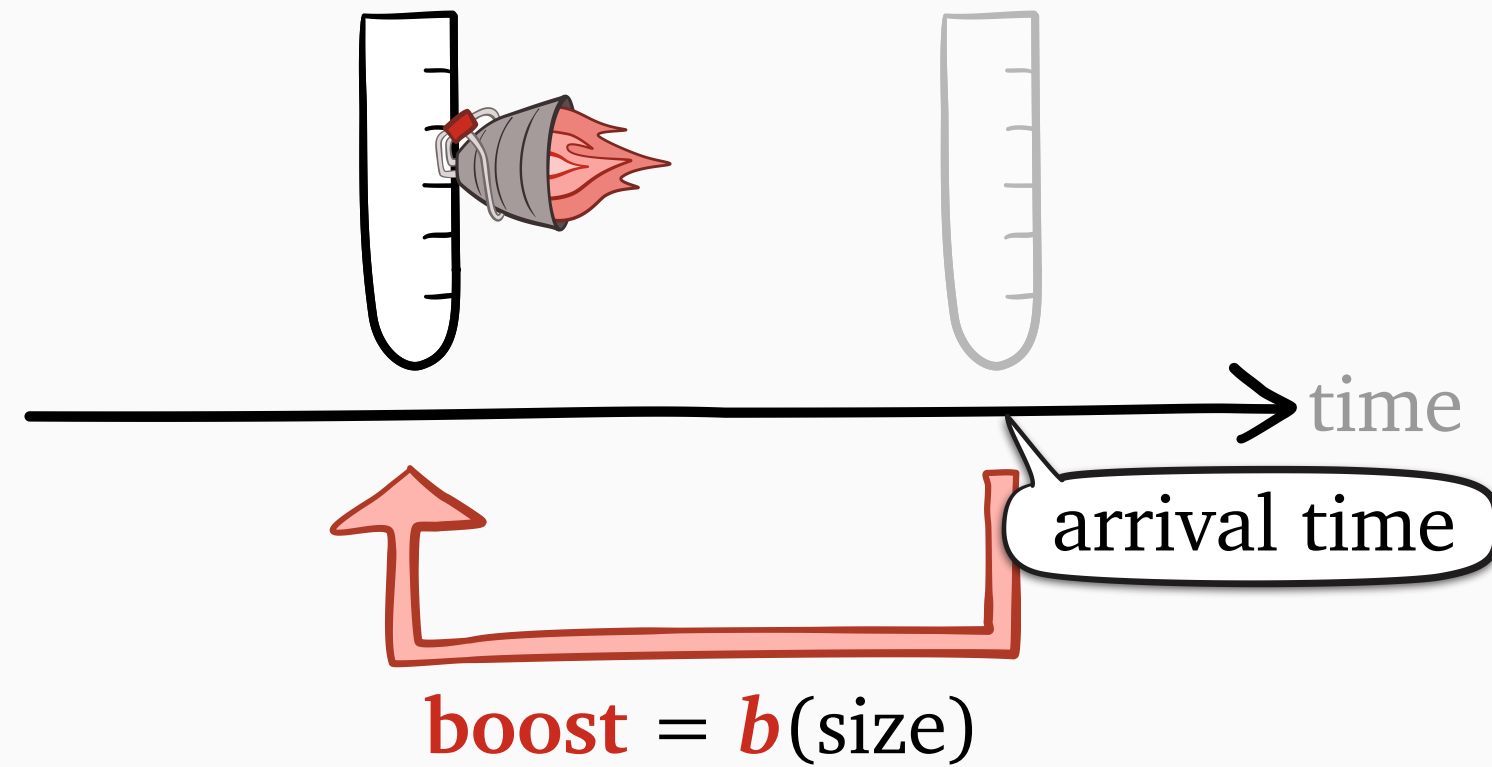
boosted arrival time
= arrival time - ***b***(size)

How **Boost** works



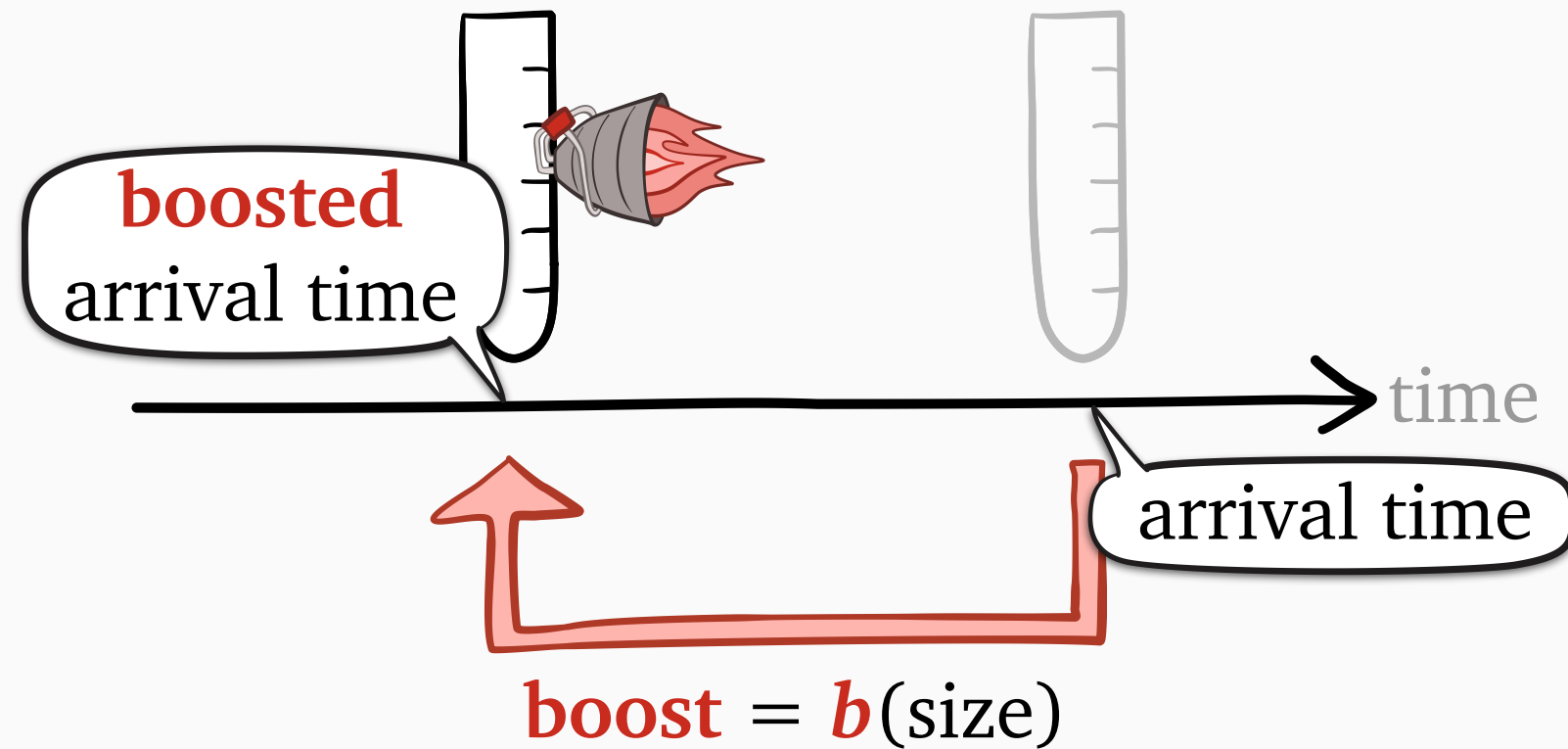
boosted arrival time
= arrival time – ***b***(size)

How **Boost** works



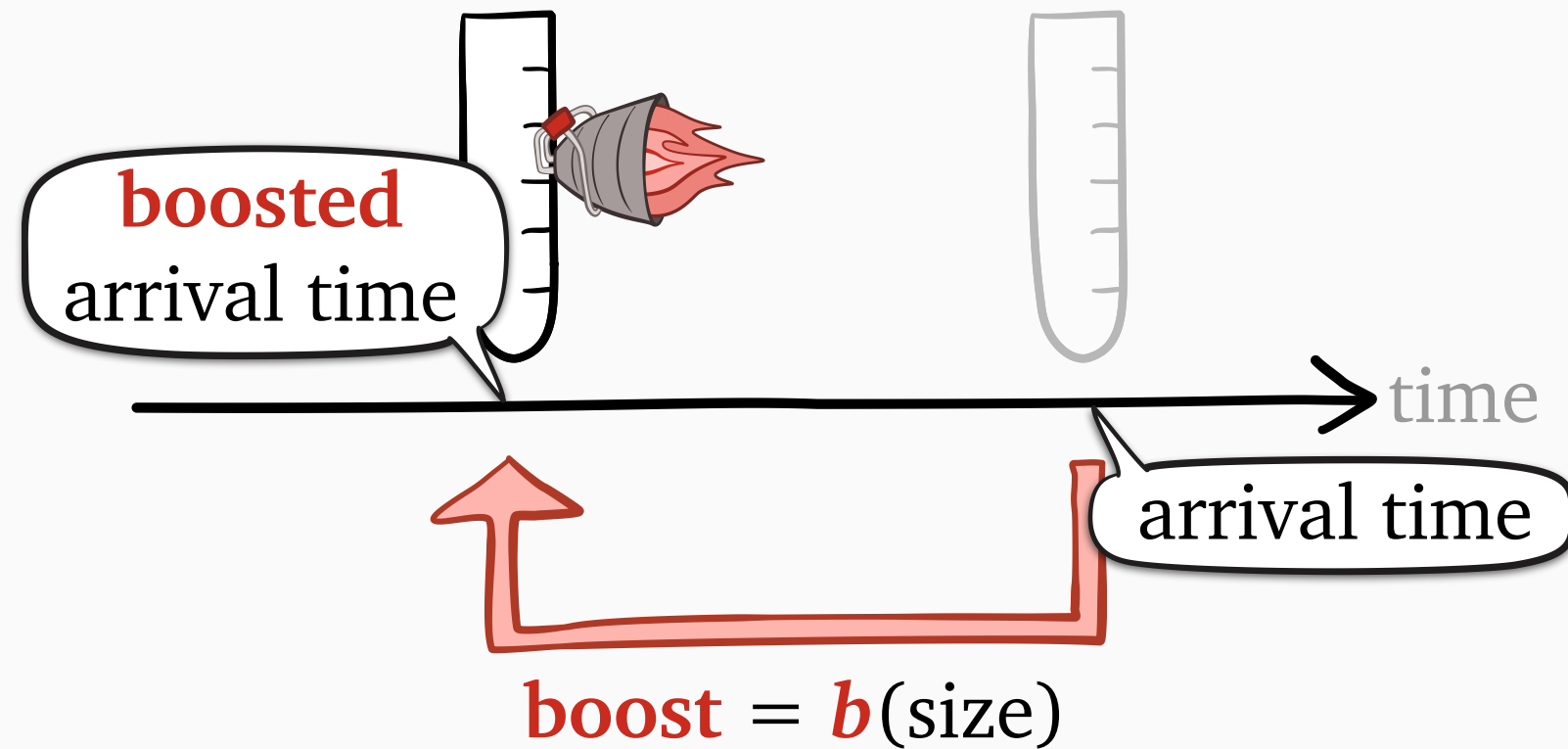
boosted arrival time
= arrival time - $b(\text{size})$

How **Boost** works

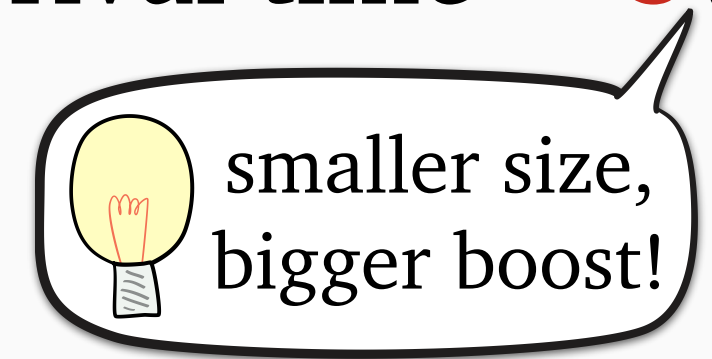


$$\text{boosted arrival time} = \text{arrival time} - b(\text{size})$$

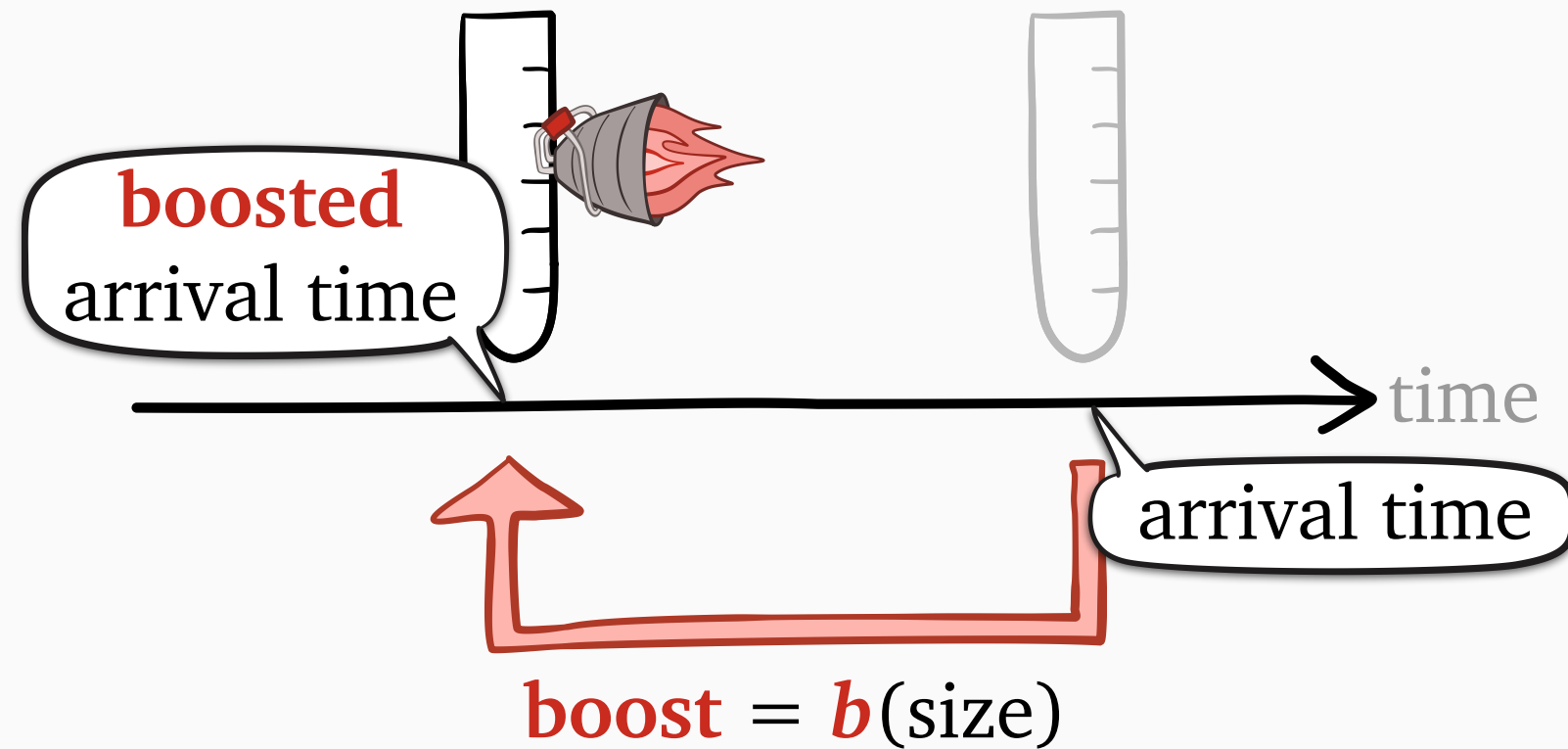
How **Boost** works



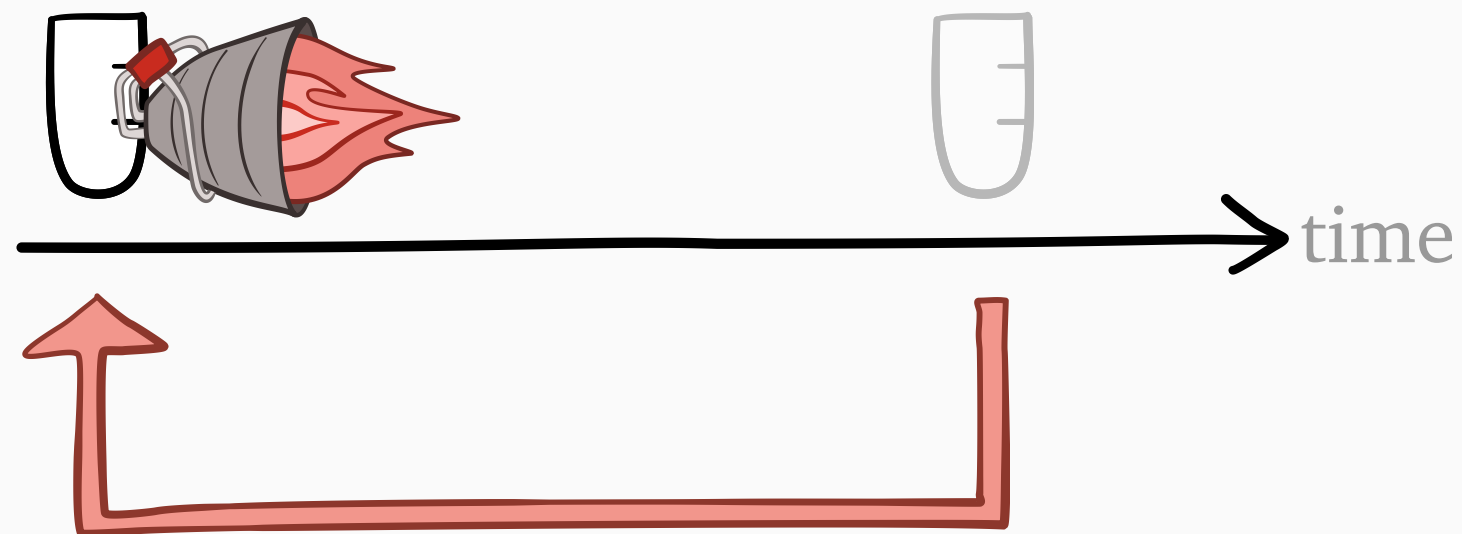
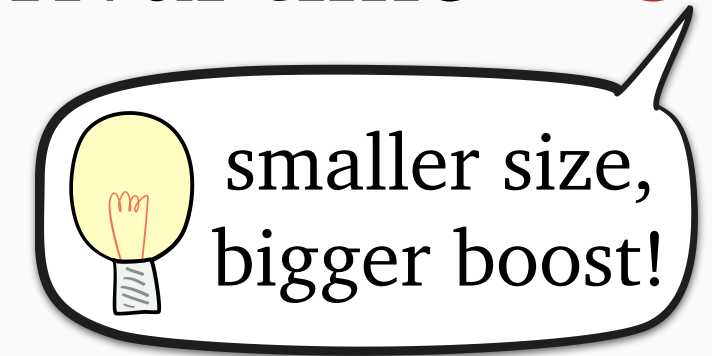
$$\text{boosted arrival time} = \text{arrival time} - b(\text{size})$$



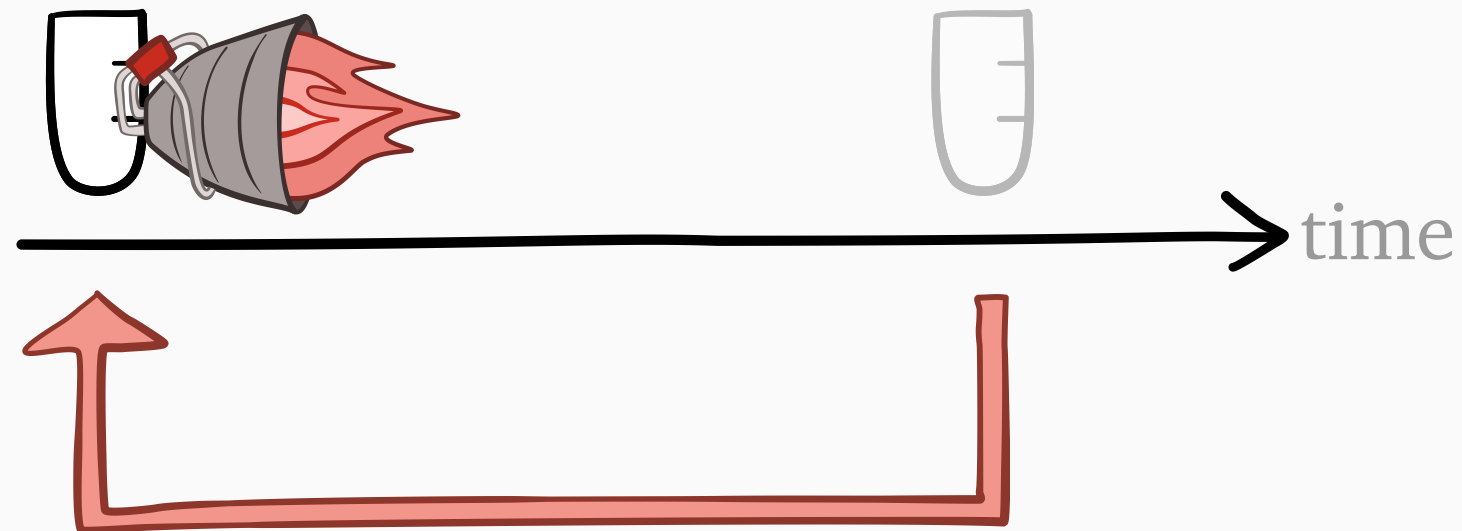
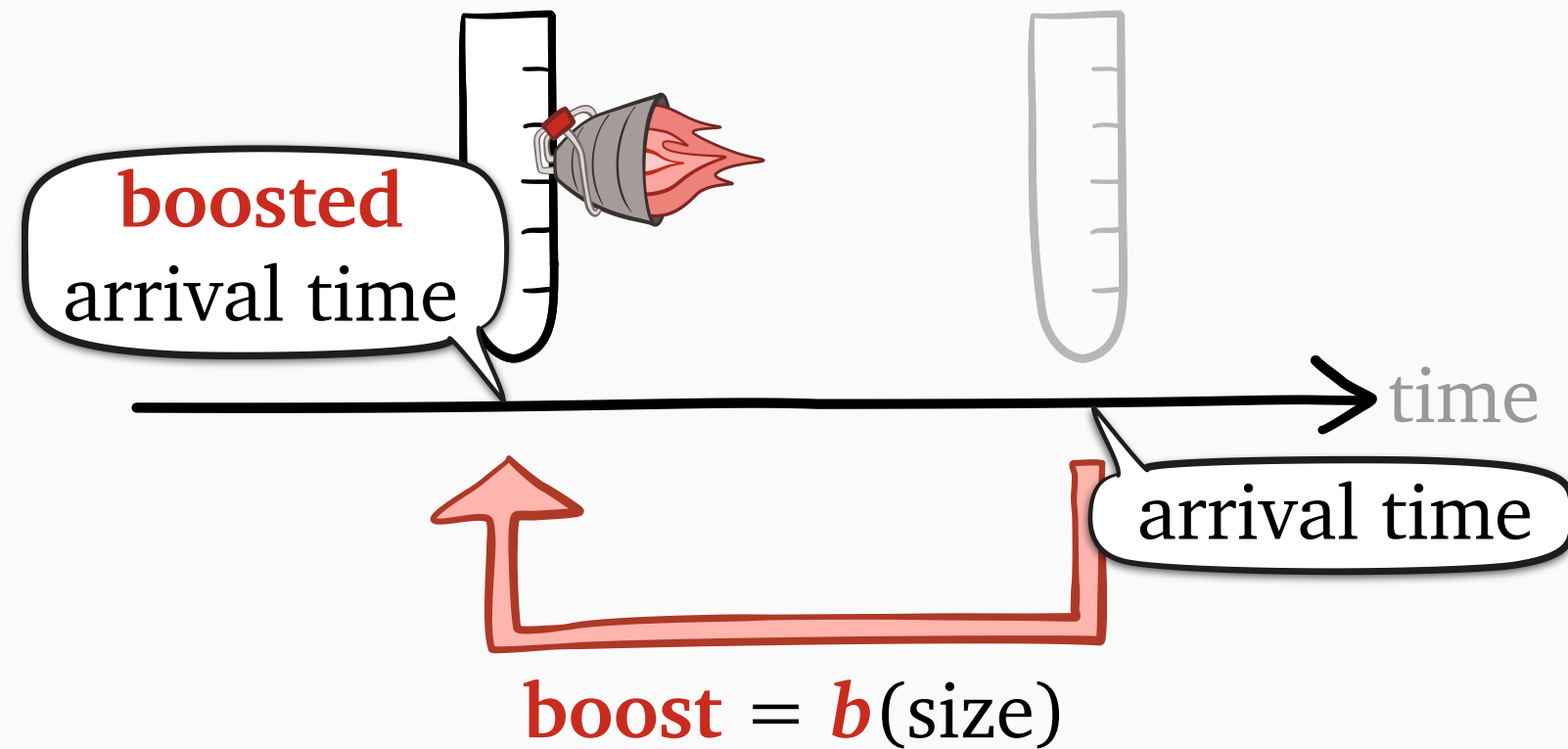
How **Boost** works



$$\text{boosted arrival time} = \text{arrival time} - \textcolor{red}{b}(\text{size})$$

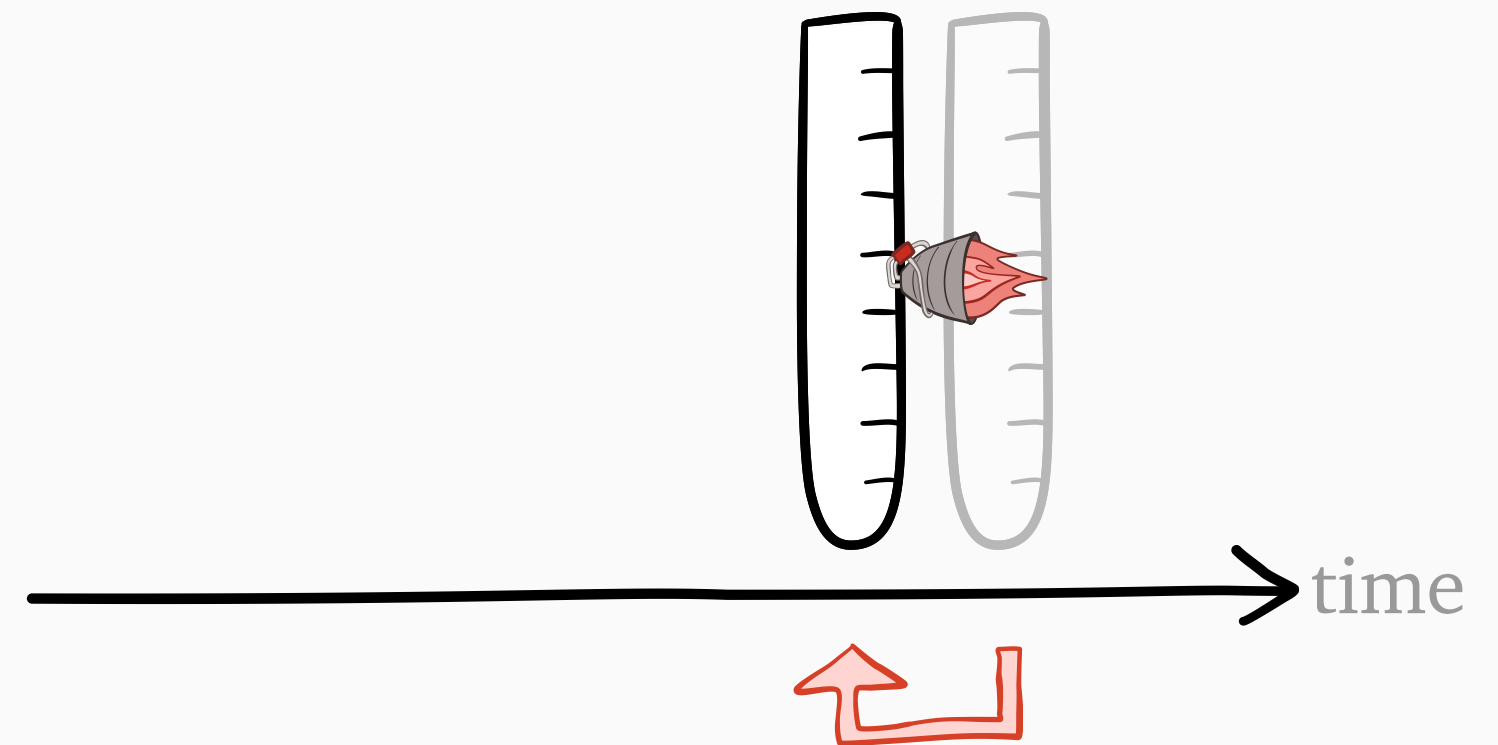


How **Boost** works

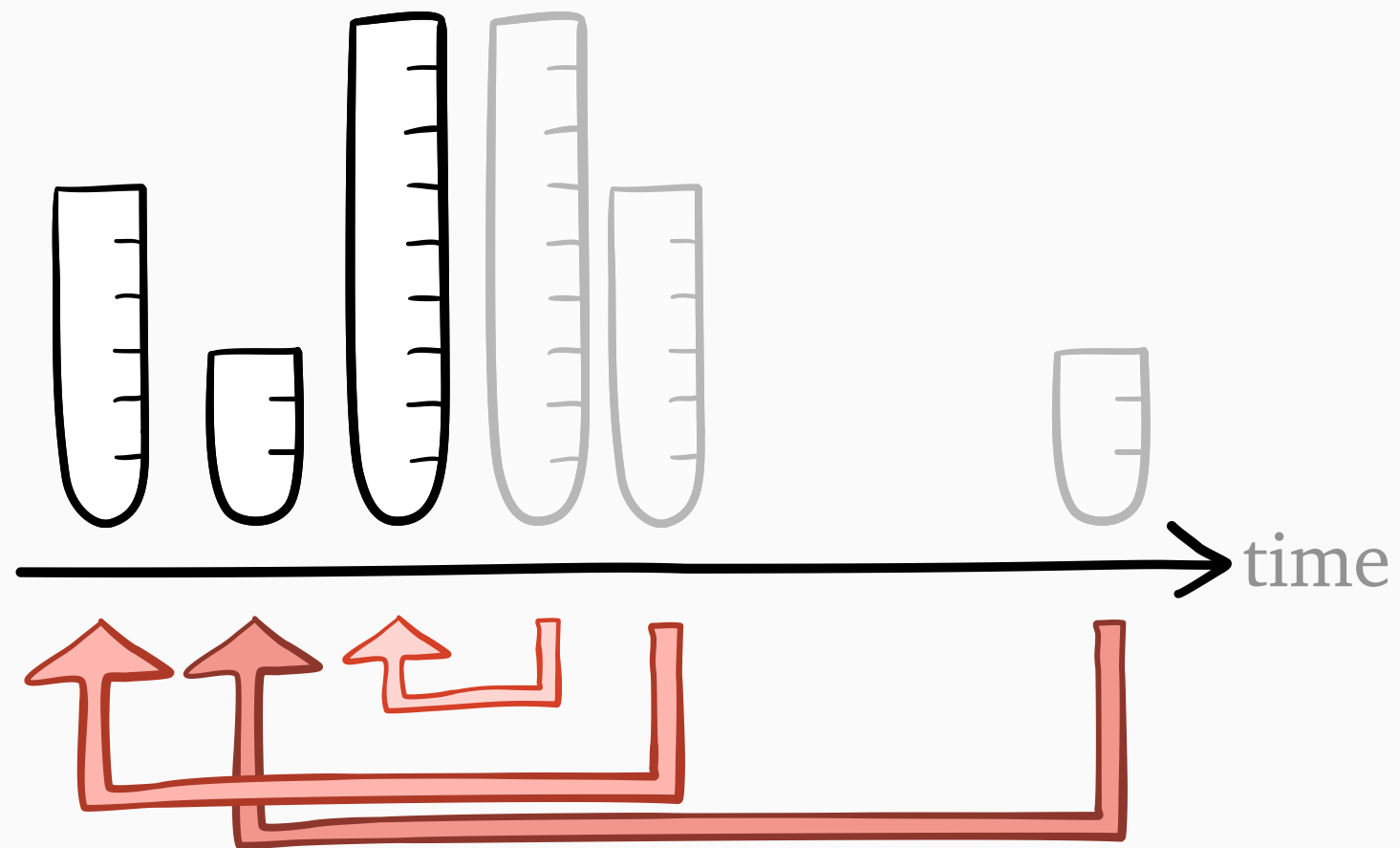


boosted arrival time
= arrival time - b (size)

💡 smaller size,
bigger boost!

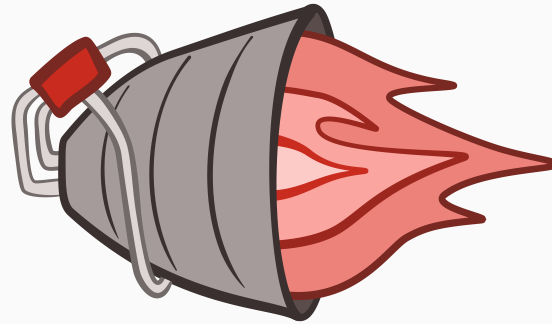


How **Boost** works

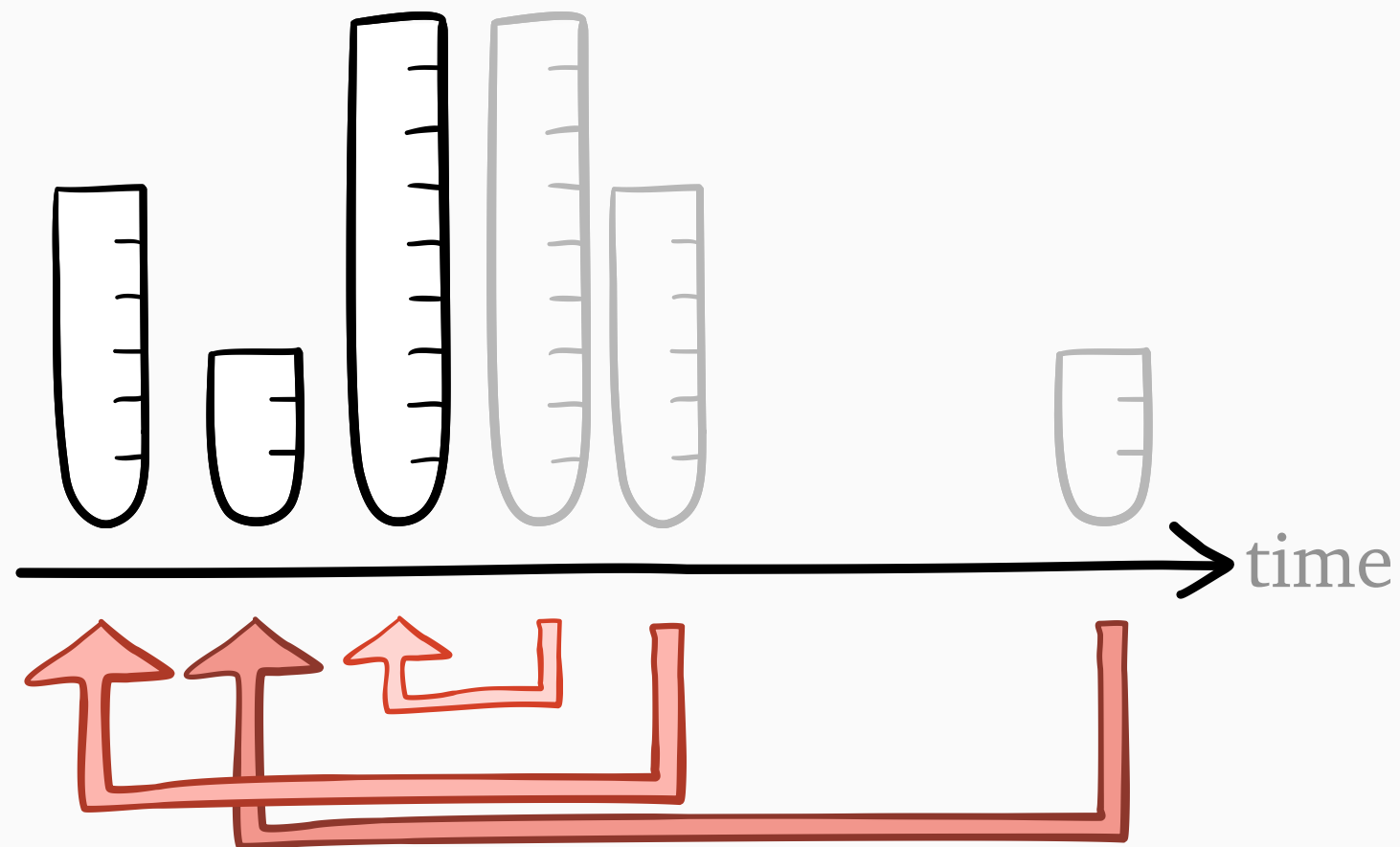


boosted arrival time
= arrival time – ***b***(size)

How **Boost** works



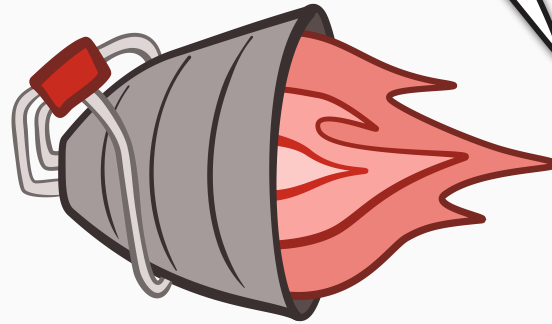
Scheduling rule: always serve job of *minimum **boosted** arrival time*



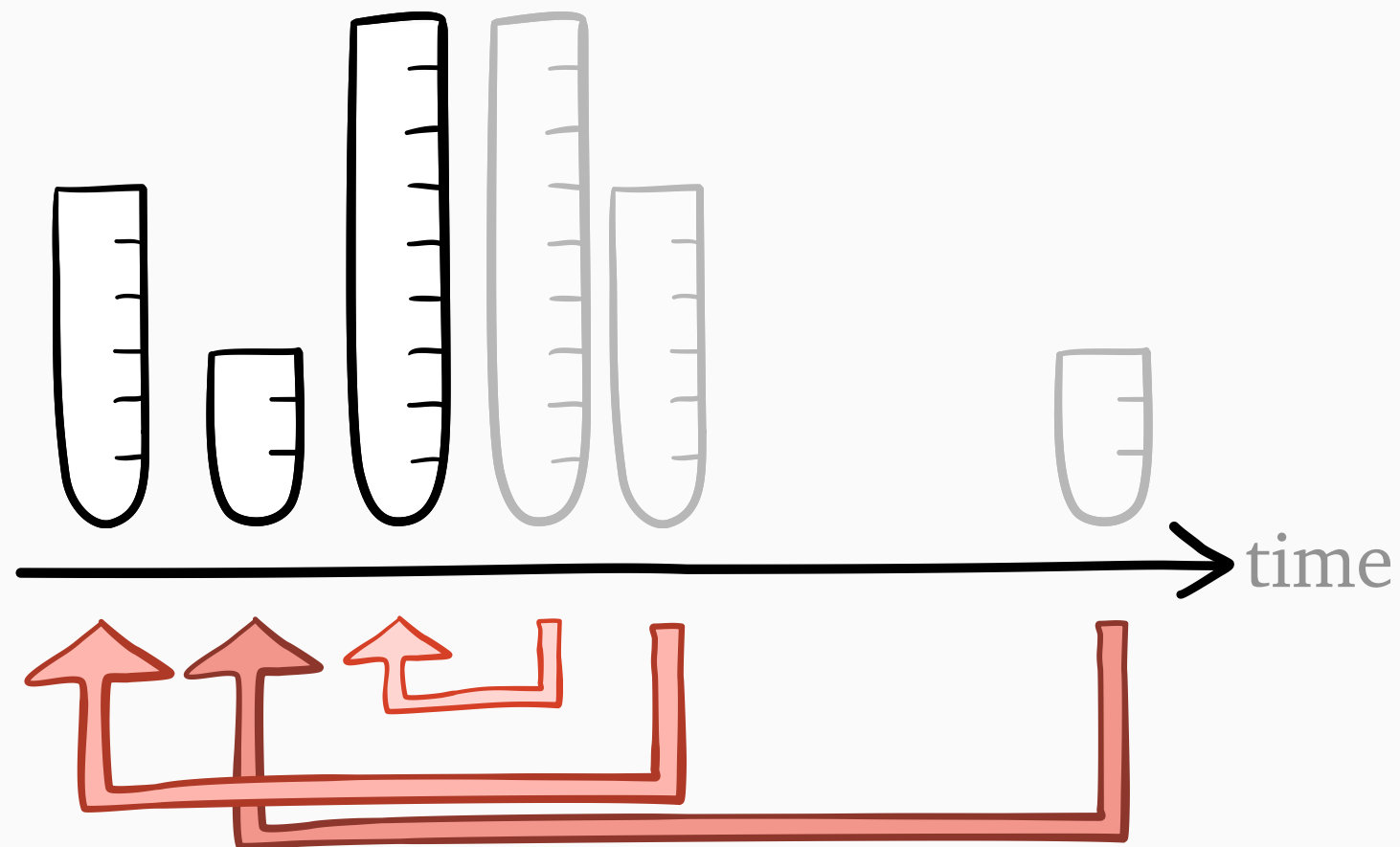
boosted arrival time
= arrival time – ***b***(size)

How **Boost** works

can be preemptive
or nonpreemptive



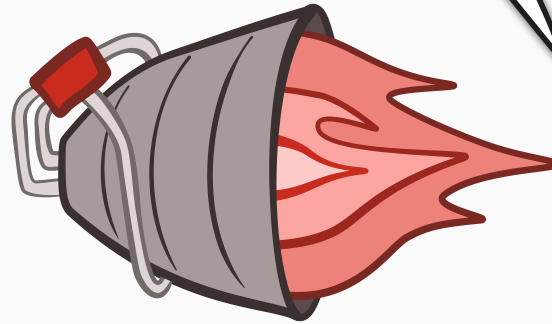
Scheduling rule: always serve job of
*minimum **boosted** arrival time*



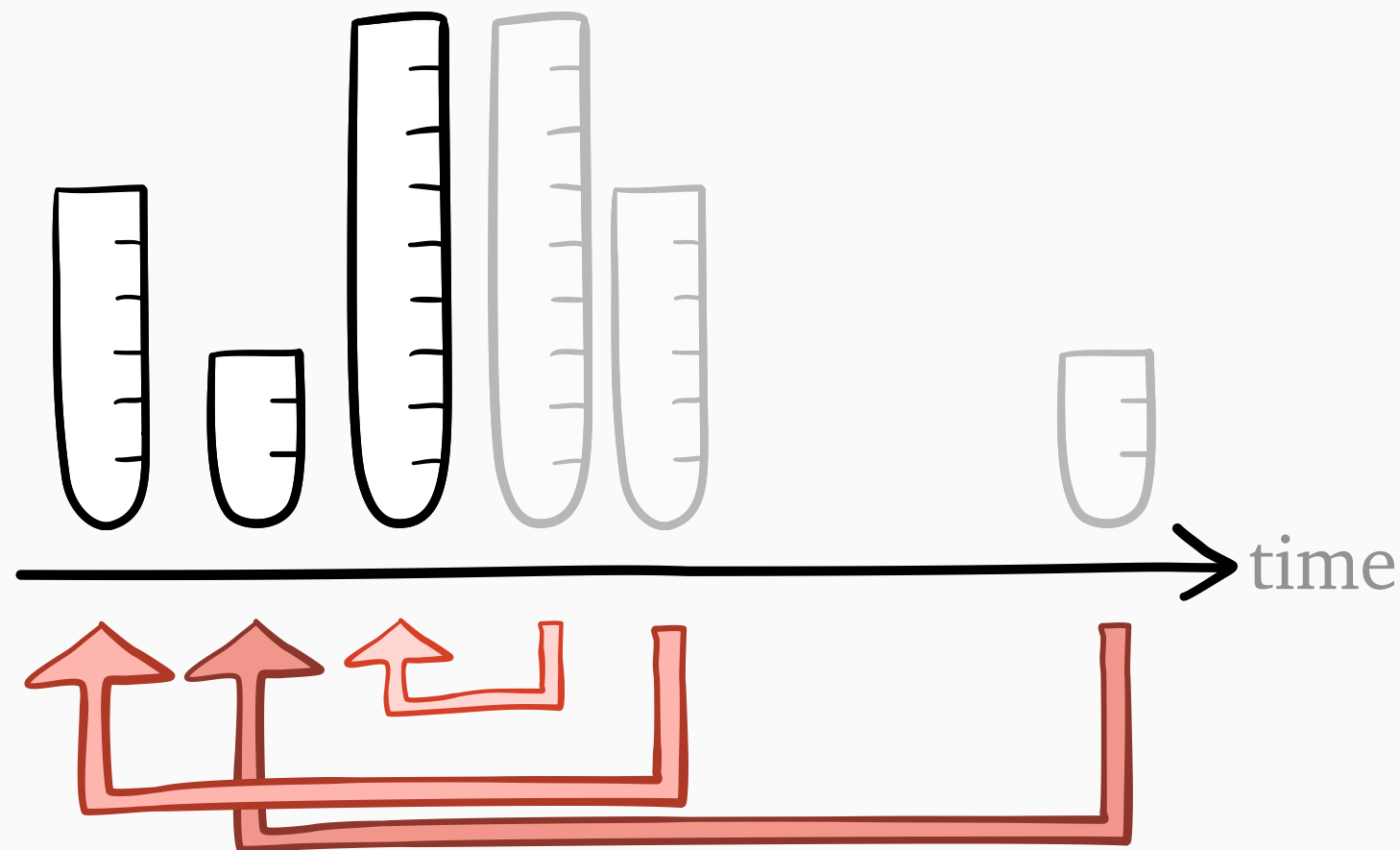
boosted arrival time
= arrival time – ***b***(size)

How **Boost** works

can be preemptive
or nonpreemptive



Scheduling rule: always serve job of
*minimum **boosted** arrival time*



boosted arrival time
= arrival time – ***b***(size)

? What's the right
boost function?

Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]$$

Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]$$

Batch problem


$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]$$



“ $\infty \cdot \mathbf{P}[T > \infty]$ ”

Batch problem


$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]$$



“ $\infty \cdot \mathbf{P}[T > \infty]$ ”

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty]$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} \underbrace{e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}}$$

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty)$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} \underbrace{e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}}$$

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

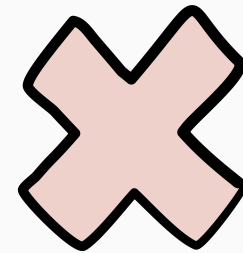
Queueing problem

$$\text{minimize } C = \lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]$$

$\underbrace{\hspace{10em}}$
“ $\infty \cdot \mathbf{P}[T > \infty]$ ”

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0$$



$$t_i = d_i - a_i$$

a_i = arrival time of job i

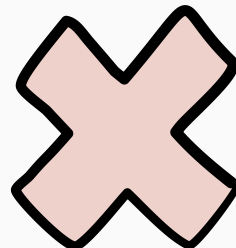
d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}} = \lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0$$
 

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}} = \underbrace{\lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]}_{\text{“}0 \cdot \mathbf{E}[e^{\gamma T}]\text{”}}$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0 \quad \times$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}} = \underbrace{\lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]}_{\text{“}0 \cdot \mathbf{E}[e^{\gamma T}]\text{”}}$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0 \quad \times$$

$$\text{minimize } \mathbf{E}[e^{\gamma T}]$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}} = \underbrace{\lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]}_{\text{“}0 \cdot \mathbf{E}[e^{\gamma T}]\text{”}}$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0 \quad \times$$

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i}$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

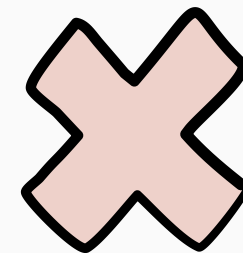
Queueing problem

$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}} = \underbrace{\lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]}_{\text{“}0 \cdot \mathbf{E}[e^{\gamma T}]\text{”}}$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0$$



$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i}$$



$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{“}\infty \cdot \mathbf{P}[T > \infty]\text{”}} = \underbrace{\lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]}_{\text{“}0 \cdot \mathbf{E}[e^{\gamma T}]\text{”}}$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > \infty) = 0 \quad \times$$

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i} \quad \checkmark$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Queueing problem

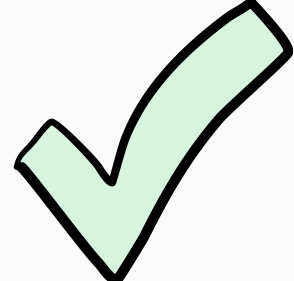
$$\text{minimize } C = \underbrace{\lim_{t \rightarrow \infty} e^{\gamma t} \mathbf{P}[T > t]}_{\text{"}\infty \cdot \mathbf{P}[T > \infty]\text{"}} = \underbrace{\lim_{\theta \rightarrow \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]}_{\text{"}0 \cdot \mathbf{E}[e^{\gamma T}]\text{"}}$$

(by final value theorem
for Laplace transforms)

Batch problem

$$\text{minimize } \mathbf{P}[T > \infty] = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(t_i > 0)$$

almost classic problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$


$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$t_i = d_i - a_i$
 $a_i = \text{arrival time of job } i$
 $d_i = \text{departure time of job } i$

Classic metric: mean weighted discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i



Classic metric: mean weighted
discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$t_i = d_i - a_i$
 $a_i = \text{arrival time of job } i$
 $d_i = \text{departure time of job } i$

$\gamma > 0$



Classic metric: mean weighted discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$\gamma > 0$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

negative
discount rate



Classic metric: mean weighted
discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$\gamma > 0$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

negative
discount rate



can't start i
before a_i

Classic metric: mean weighted
discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$\gamma > 0$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

negative
discount rate



Classic metric: mean weighted
discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

can't start i
before a_i

Relaxation solved by (sign-flipped) WDSPT, which is **Boost** with

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$\gamma > 0$

$$t_i = d_i - a_i$$

a_i = arrival time of job i

d_i = departure time of job i

negative
discount rate



Classic metric: mean weighted
discounted departure time

$$\frac{1}{n} \sum_{i=1}^n w_i e^{-\theta d_i}$$

can't start i
before a_i

Relaxation solved by (sign-flipped) WDSPT, which is **Boost** with

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

γ -Boost

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$\gamma > 0$

$$t_i = d_i - a_i$$

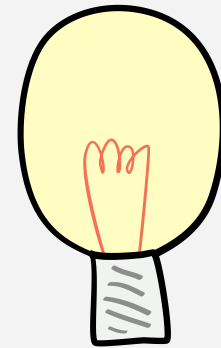
a_i = arrival time of job i

d_i = departure time of job i

negative
discount rate



can't start i
before a_i



Takeaway:

$\mathbf{E}[e^{\gamma T}]$ for critical γ is a
“smoothed” tail metric

$$\sum_{i=1}^n w_i e^{-\theta d_i}$$

Relaxation solved by (sign-flipped) WDSPT, which is **Boost** with

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

γ -Boost

Solving the batch problem

$$\text{minimize } \mathbf{E}[e^{\gamma T}] = \frac{1}{n} \sum_{i=1}^n e^{\gamma t_i} = \frac{1}{n} \sum_{i=1}^n e^{-\gamma a_i} e^{\gamma d_i}$$

$\gamma > 0$

$$t_i = d_i - a_i$$

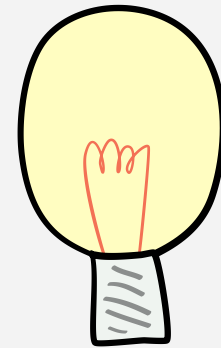
a_i = arrival time of job i

d_i = departure time of job i

negative
discount rate



can't start i
before a_i



Takeaway:

$\mathbf{E}[e^{\gamma T}]$ for critical γ is a
“smoothed” tail metric

$$\sum_{i=1}^n w_i e^{-\theta d_i}$$

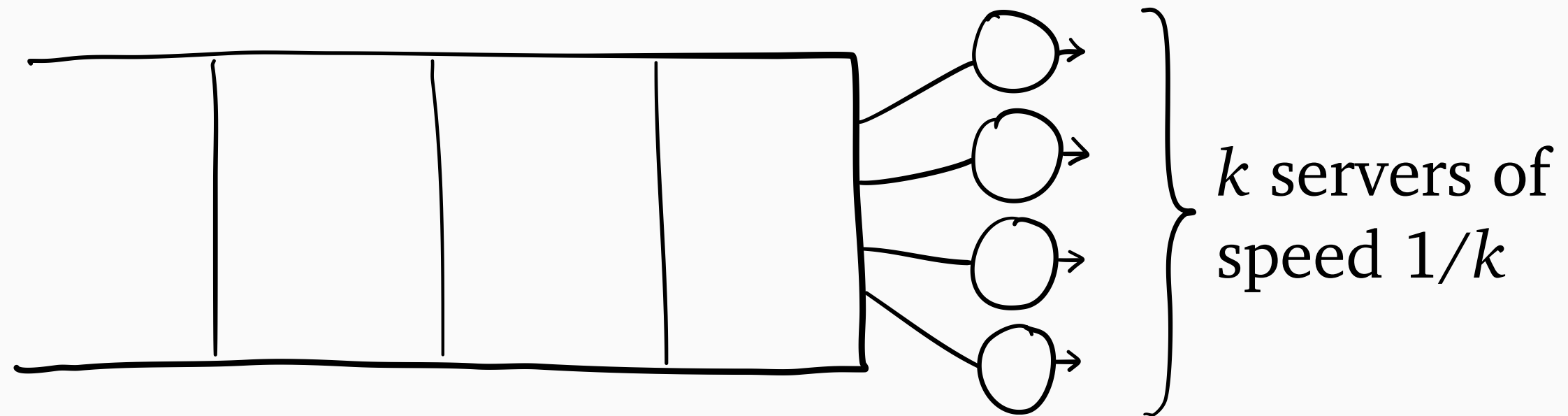
Relaxation solved by (sign-flipped) WDSPT, which is **Boost** with

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

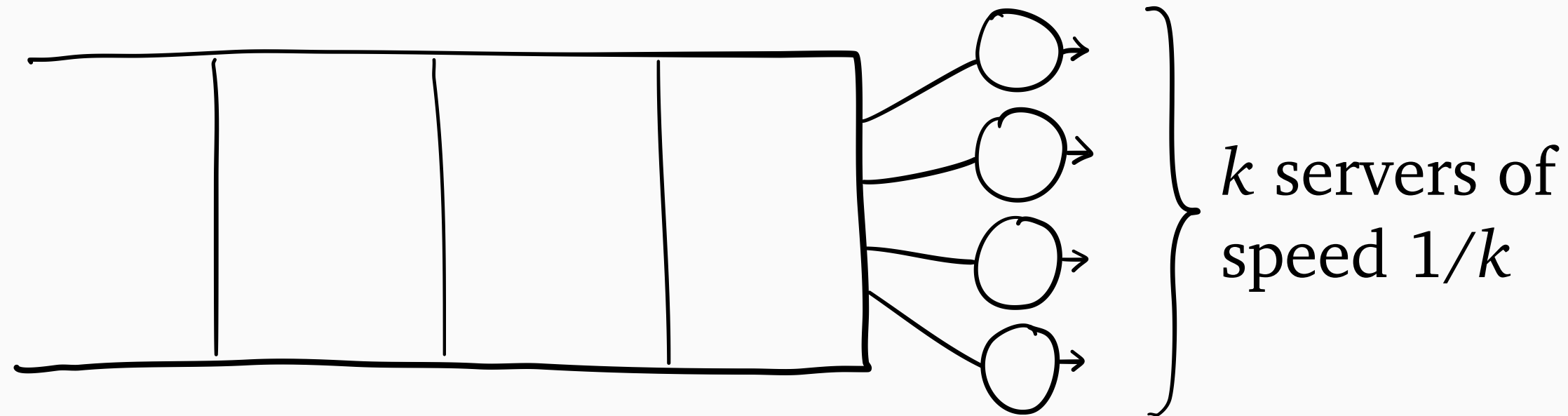
γ -Boost

Unknown sizes:
swap WDSPT for *Gittins*

Multiserver tail optimization



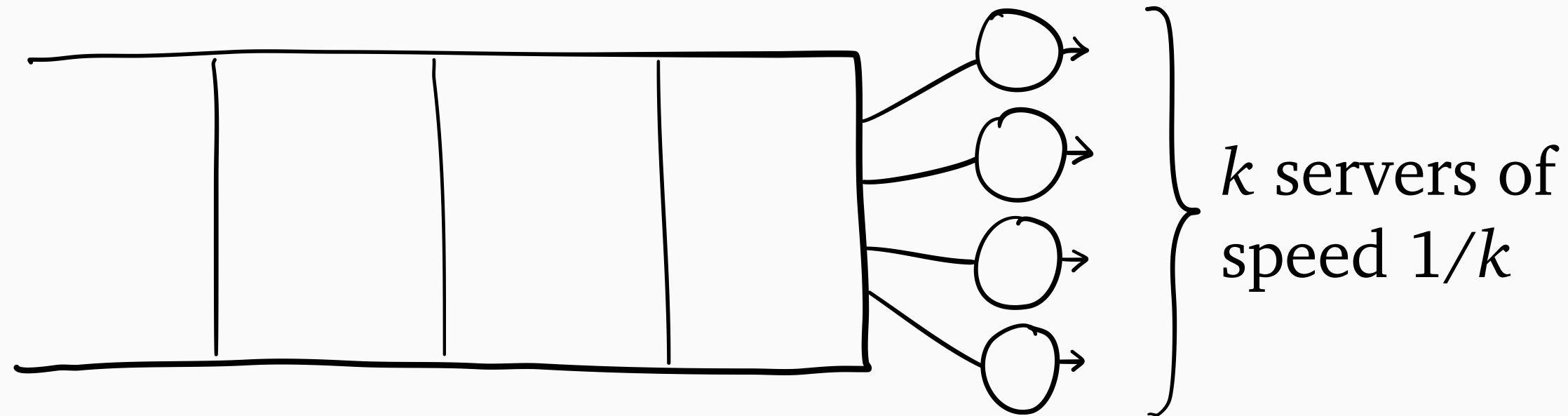
Multiserver tail optimization



Theorem: γ -Boost is strongly tail-optimal
in the heavy-traffic limit [Yu et al., 2025]

$$\rho = \lambda \mathbb{E}[S] \rightarrow 1$$

Multiserver tail optimization

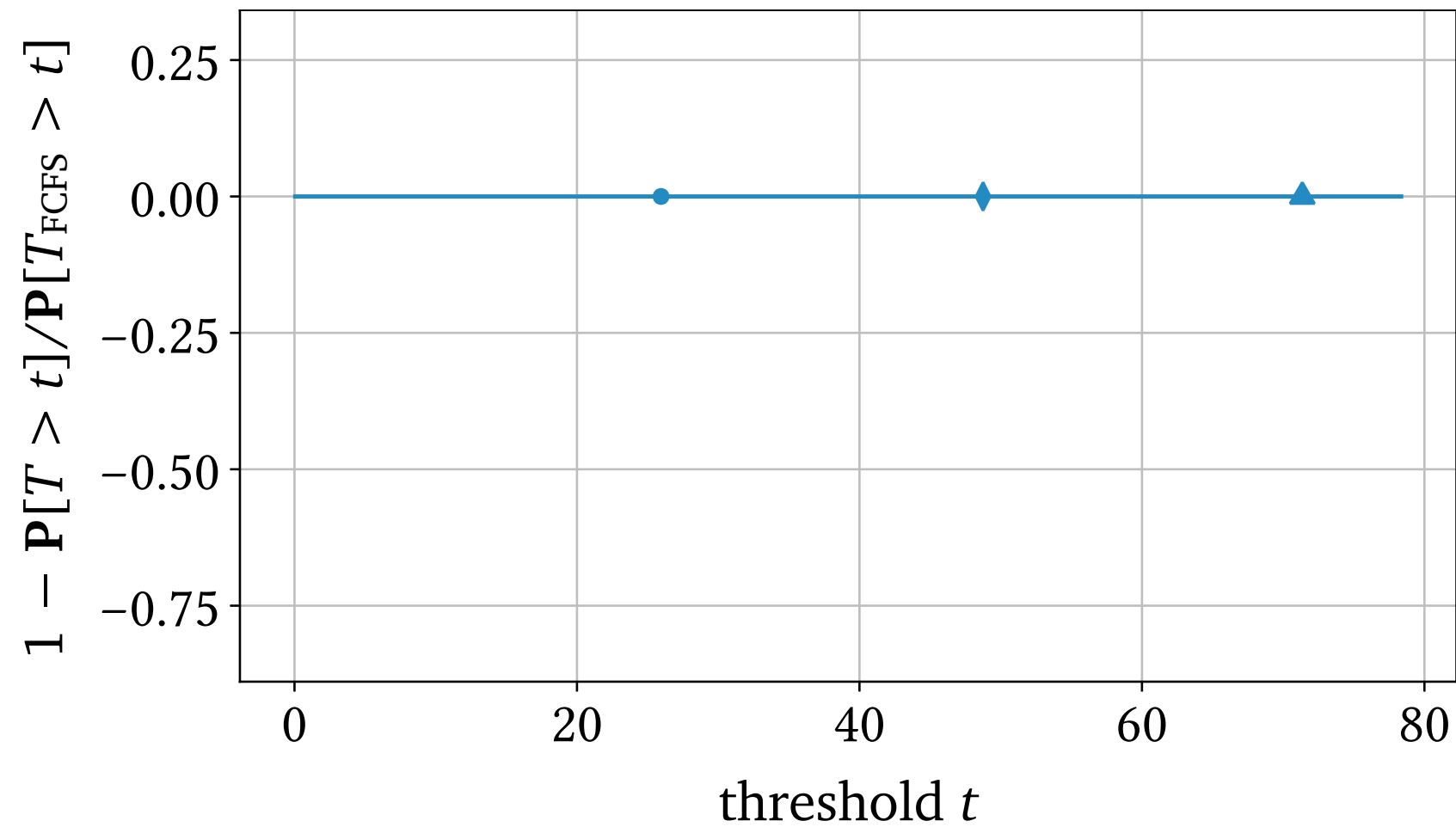


Theorem: γ -Boost is strongly tail-optimal
in the heavy-traffic limit [Yu et al., 2025]

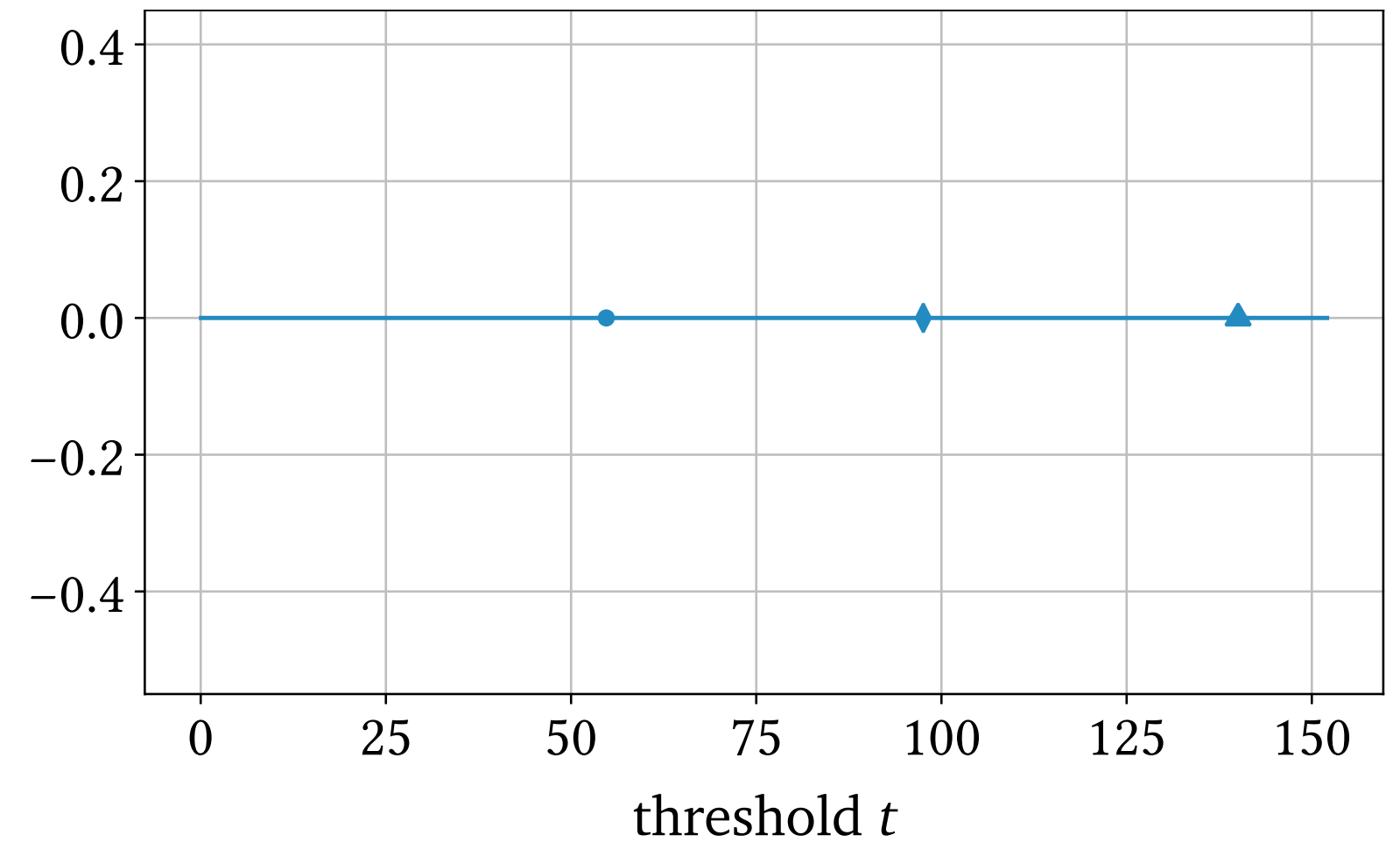
$$\rho = \lambda \mathbb{E}[S] \rightarrow 1$$

? load < 1 ?

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



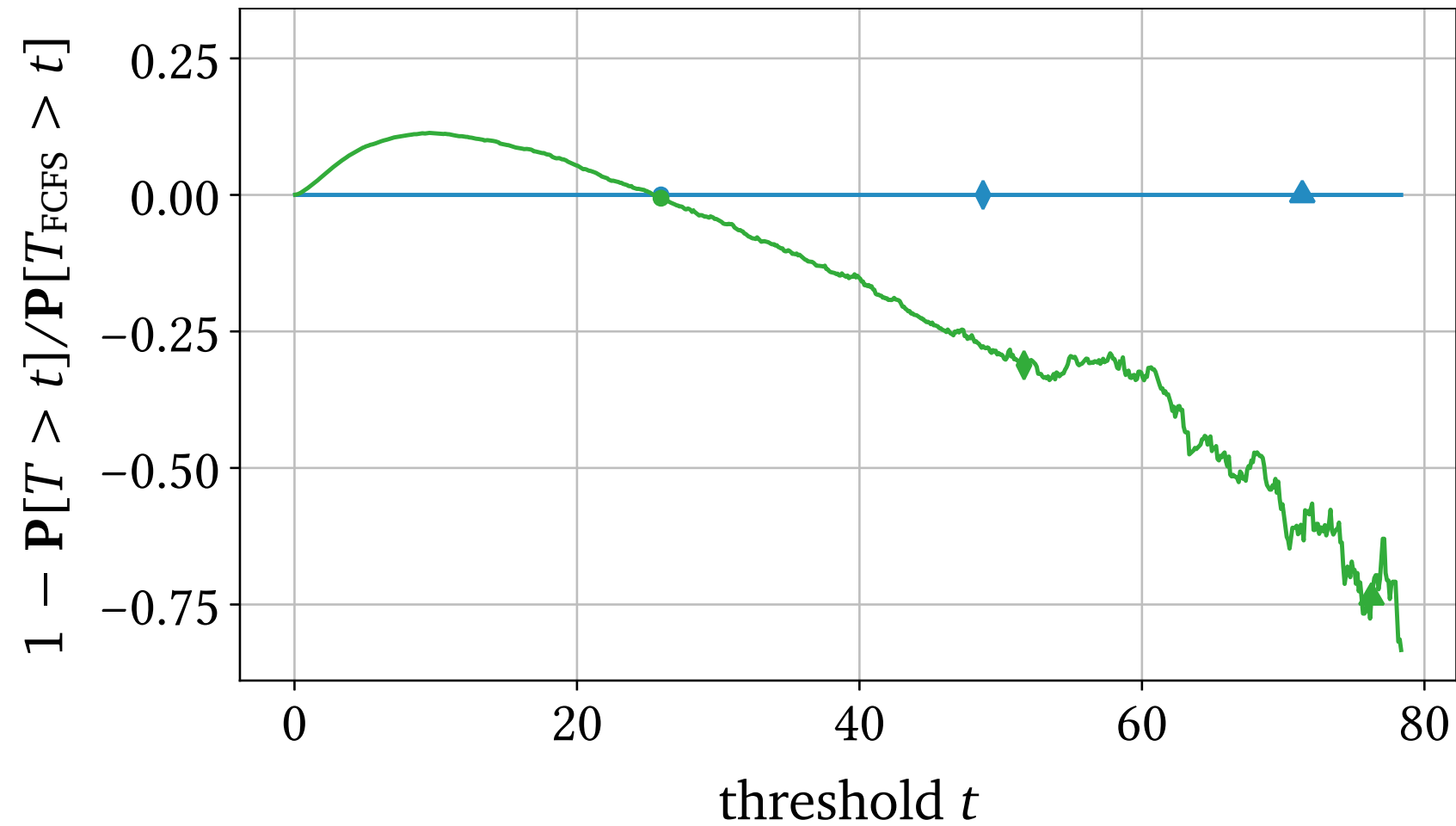
— FCFS

● 90th percentile

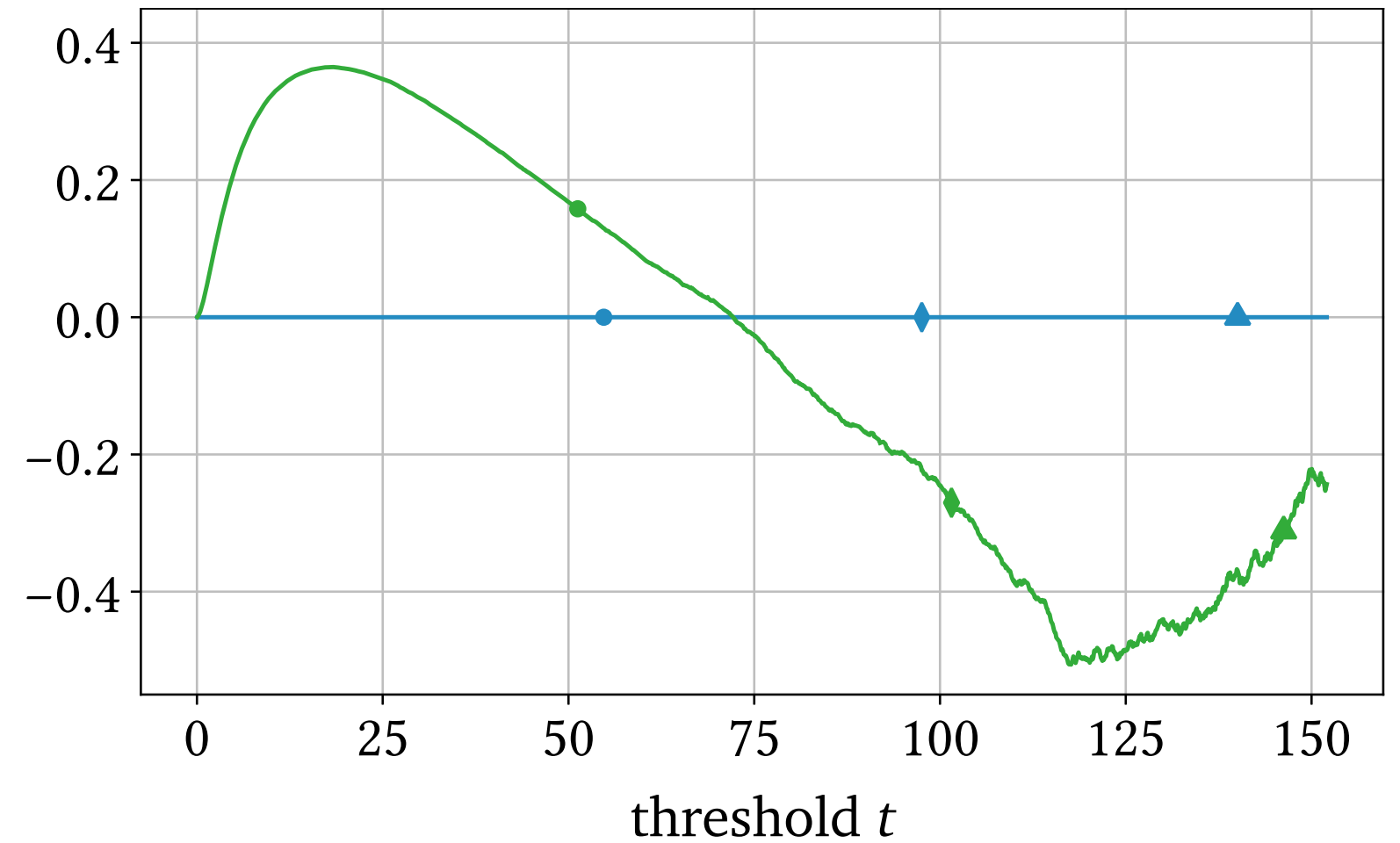
◆ 99th percentile

▲ 99.9th percentile

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



— FCFS

— γ -Boost

● 90th percentile

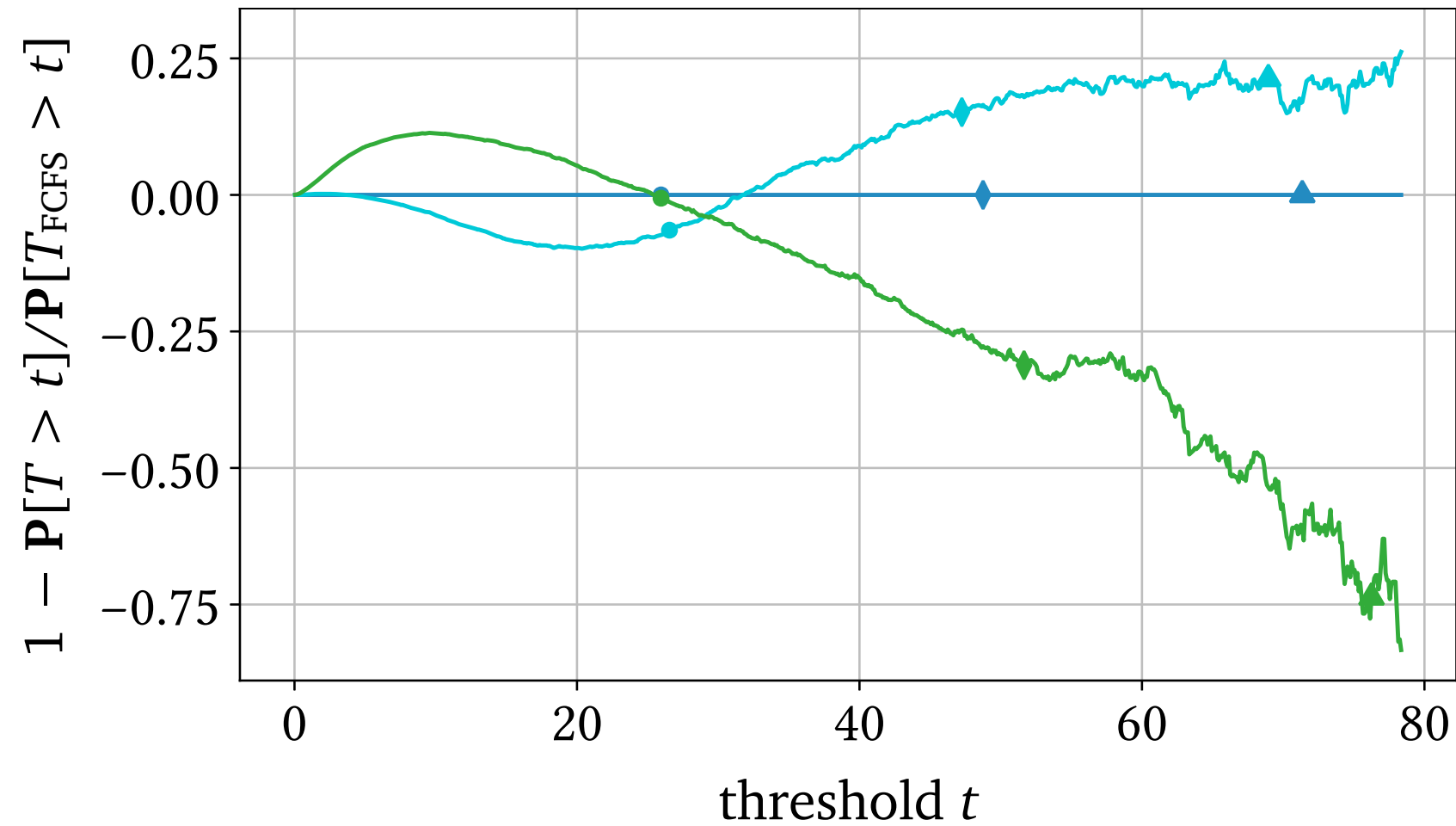
$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

● 95th percentile

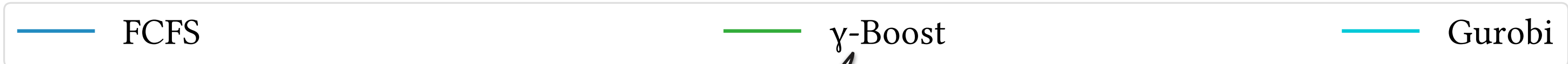
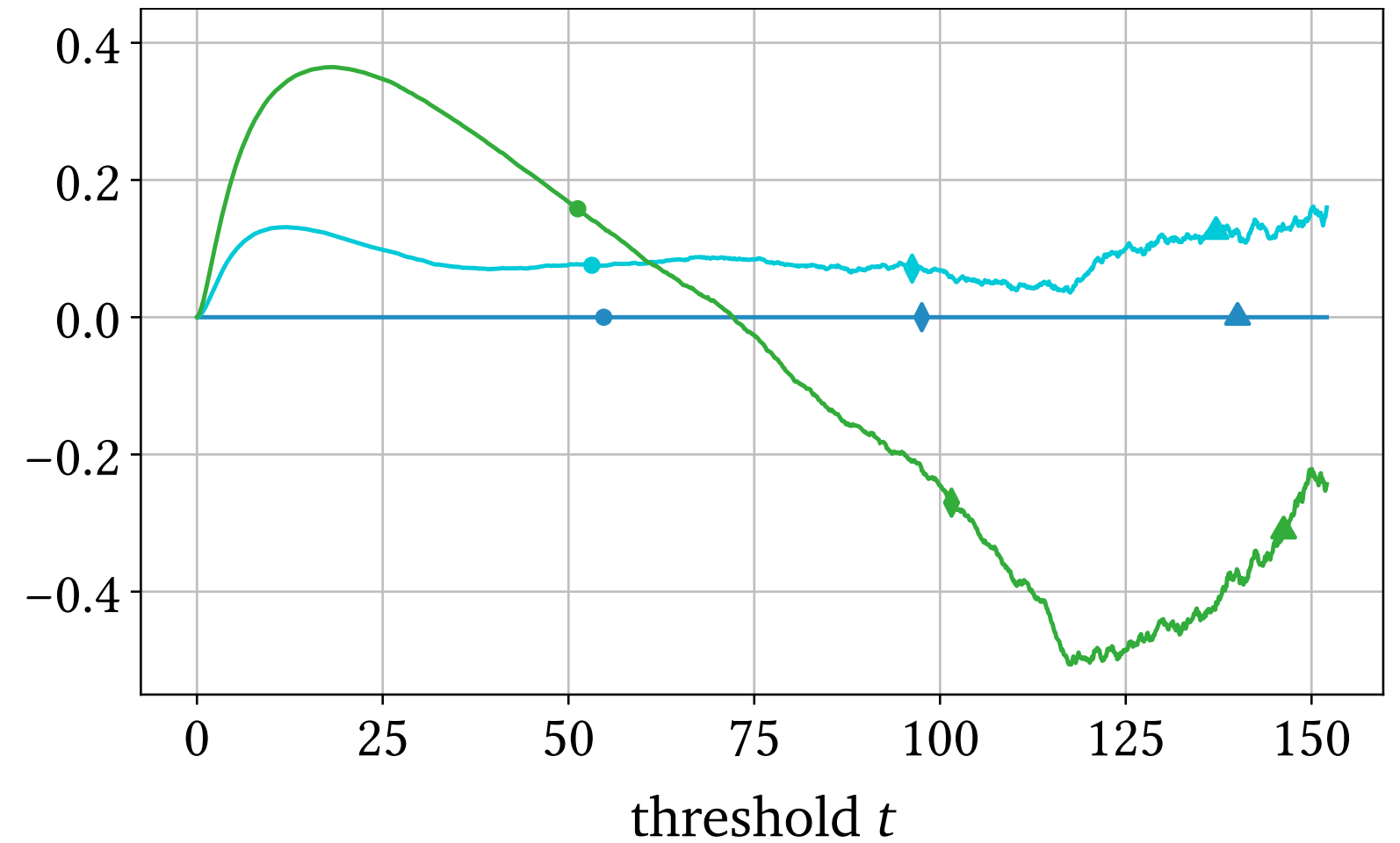
▲

99.9th percentile

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



● 90th percentile

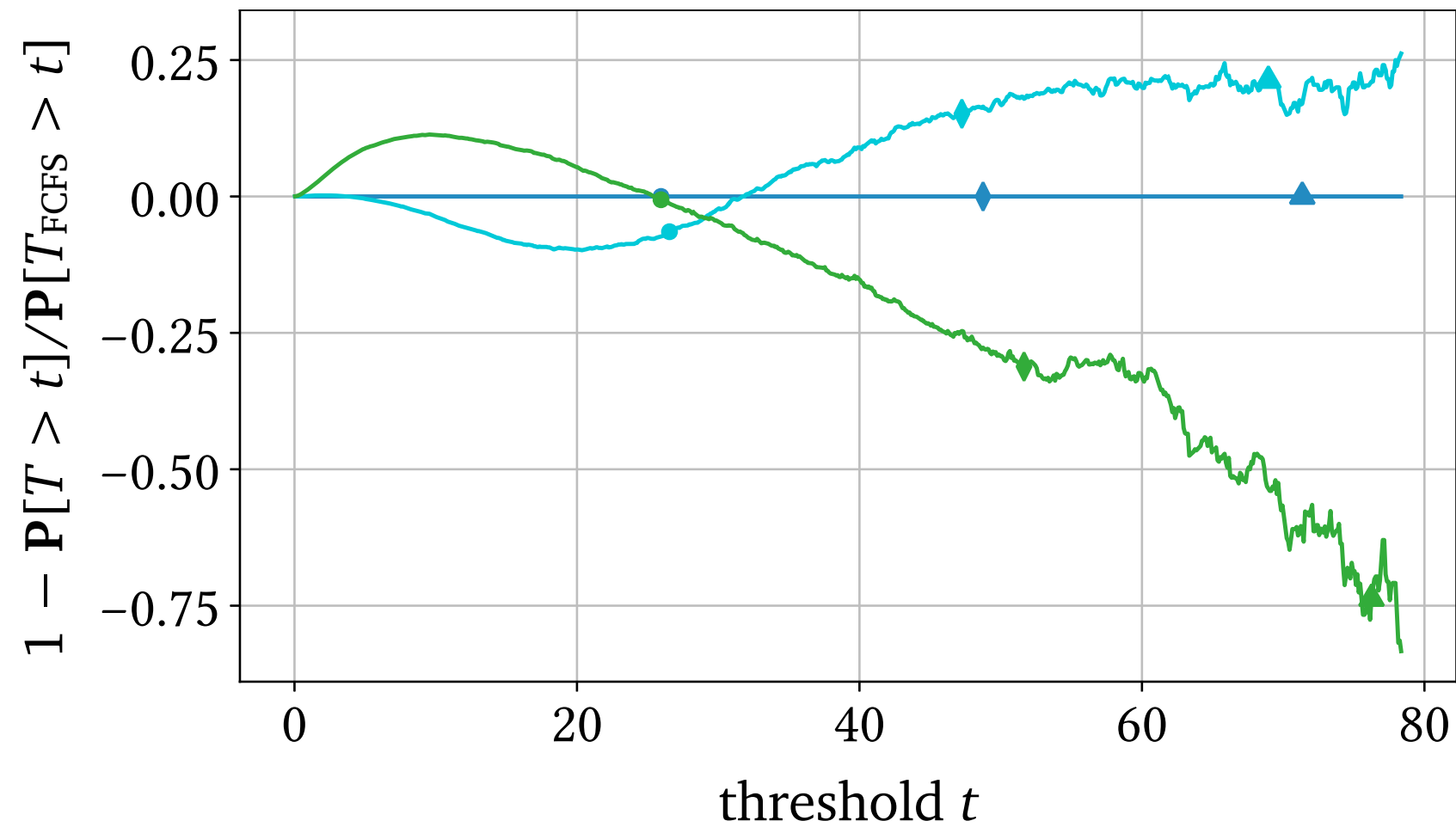
$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

● 95th percentile

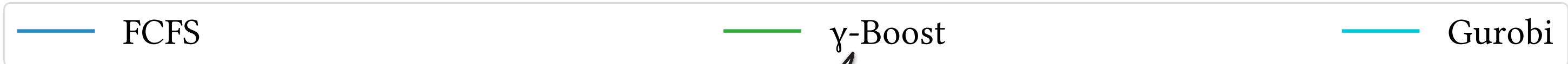
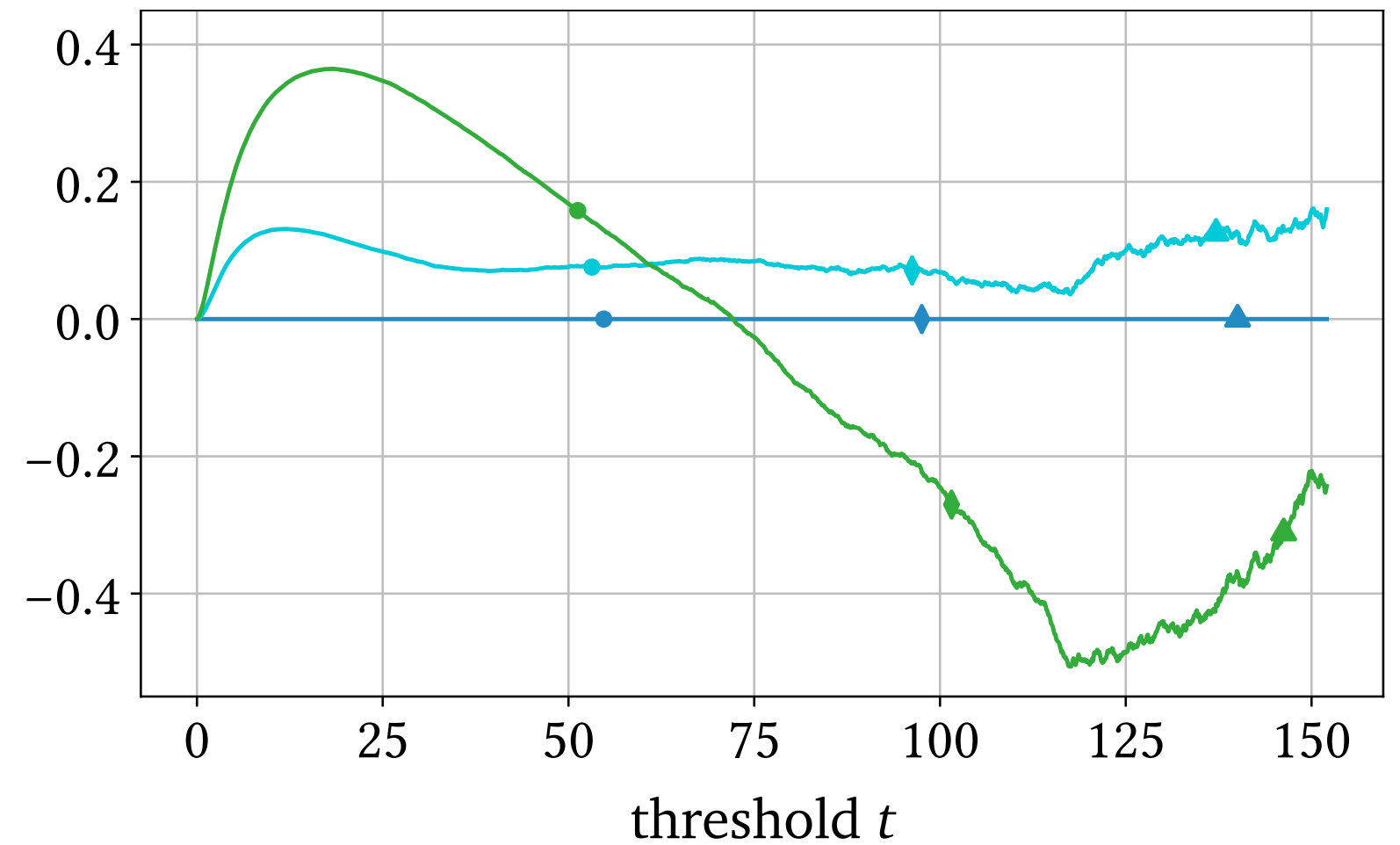
▲

99.9th percentile

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



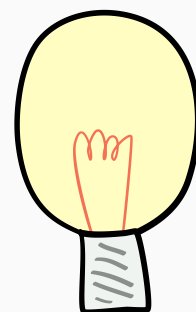
● 90th percentile

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

● 90th percentile

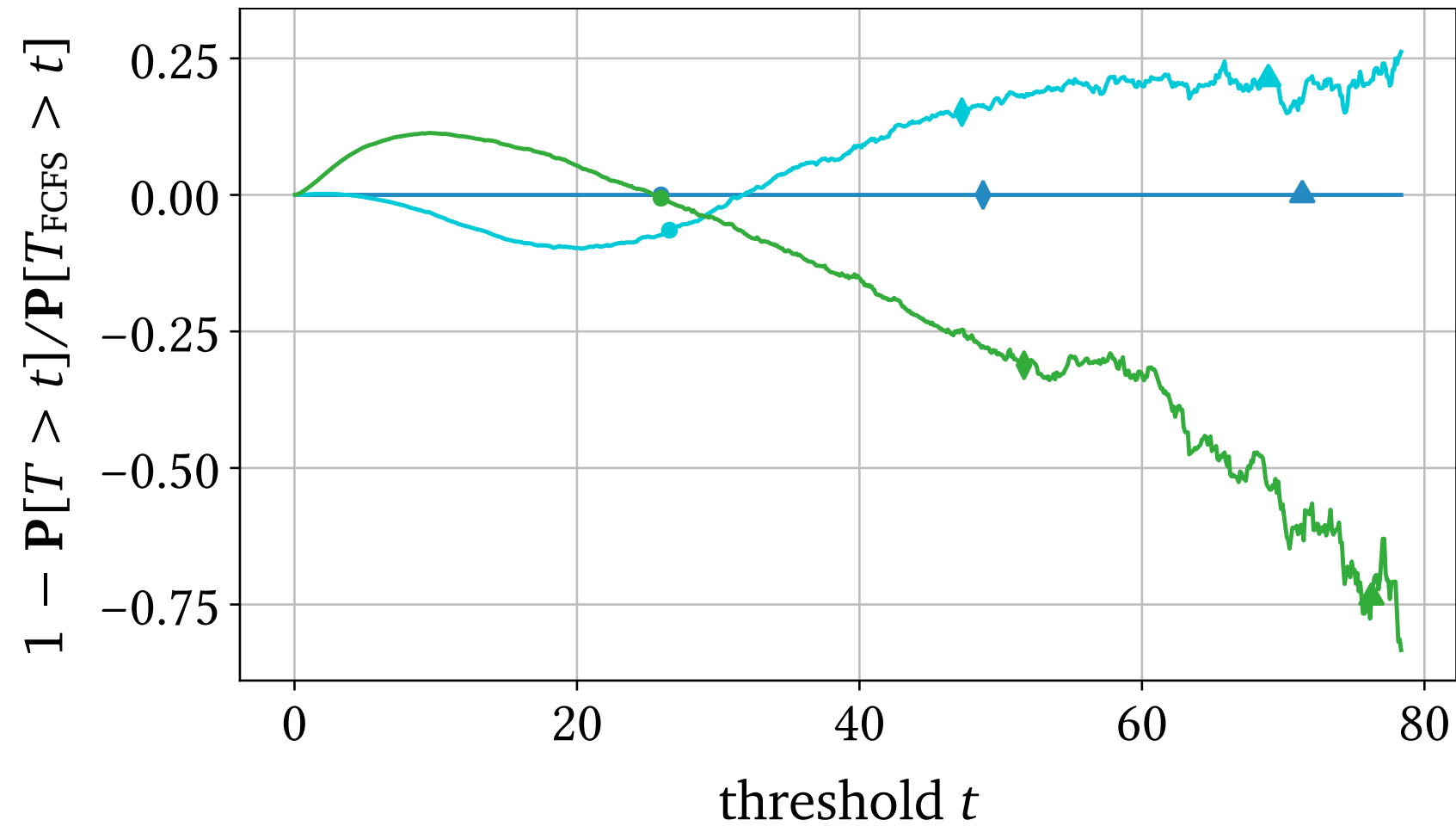


▲ 99.9th percentile

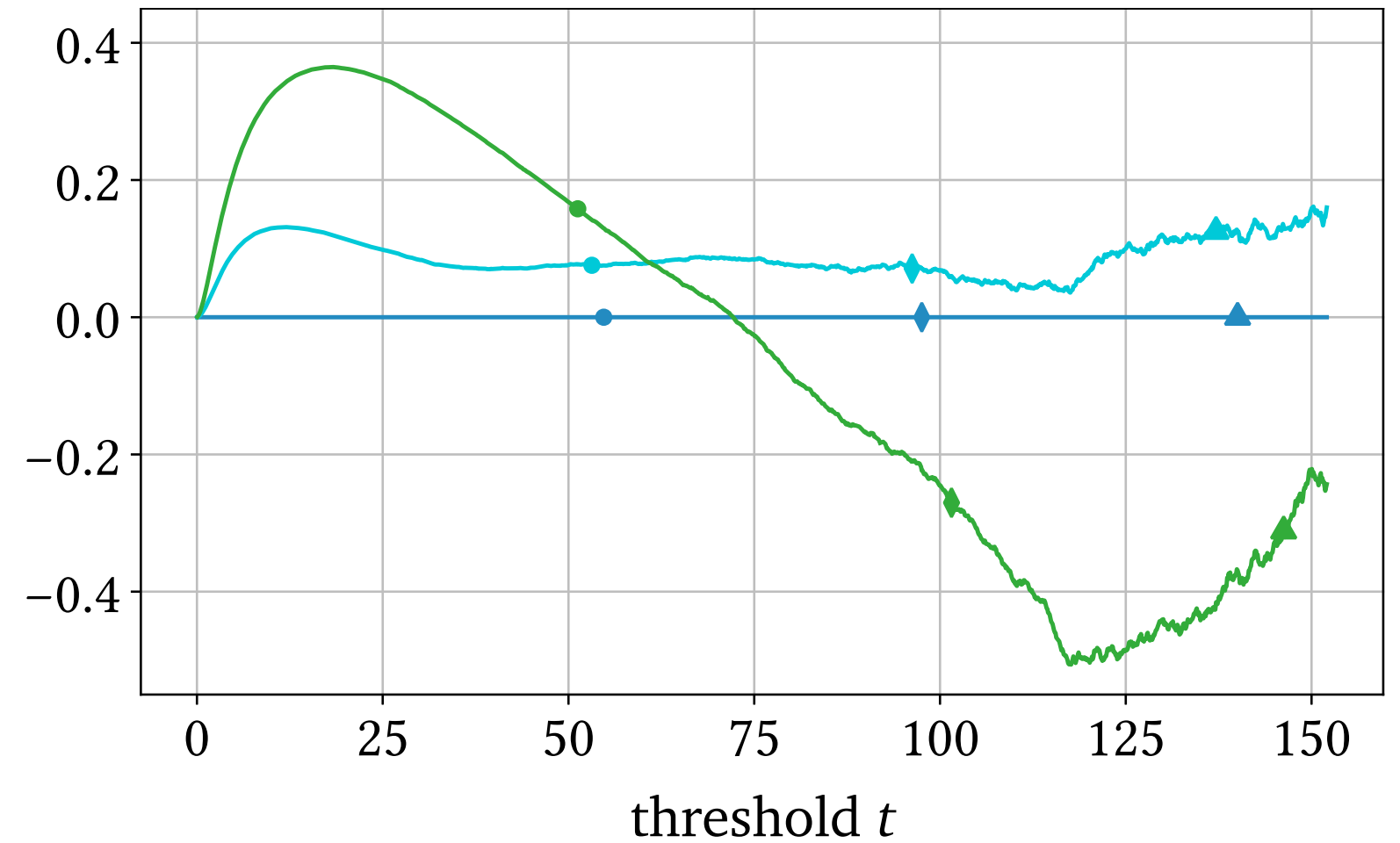


Boost large jobs?

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



● 90th percentile

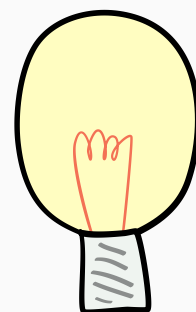
$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

● 99th percentile

▲

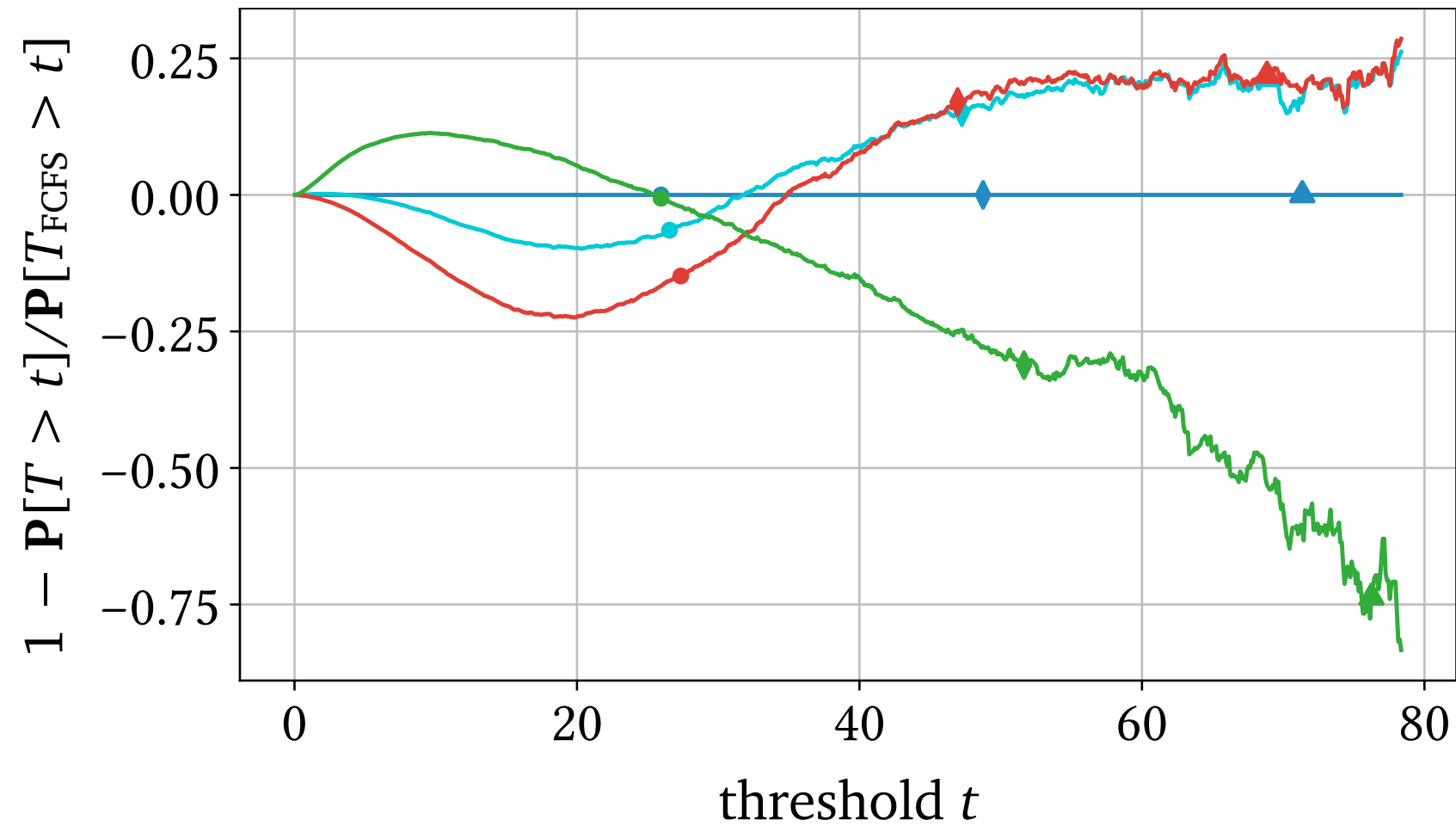
99.9th percentile

compare to bin packing

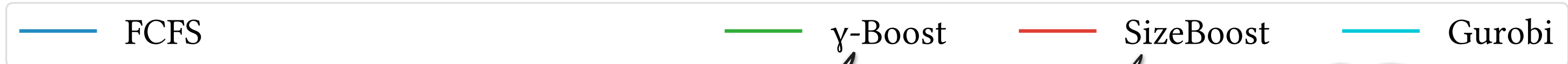
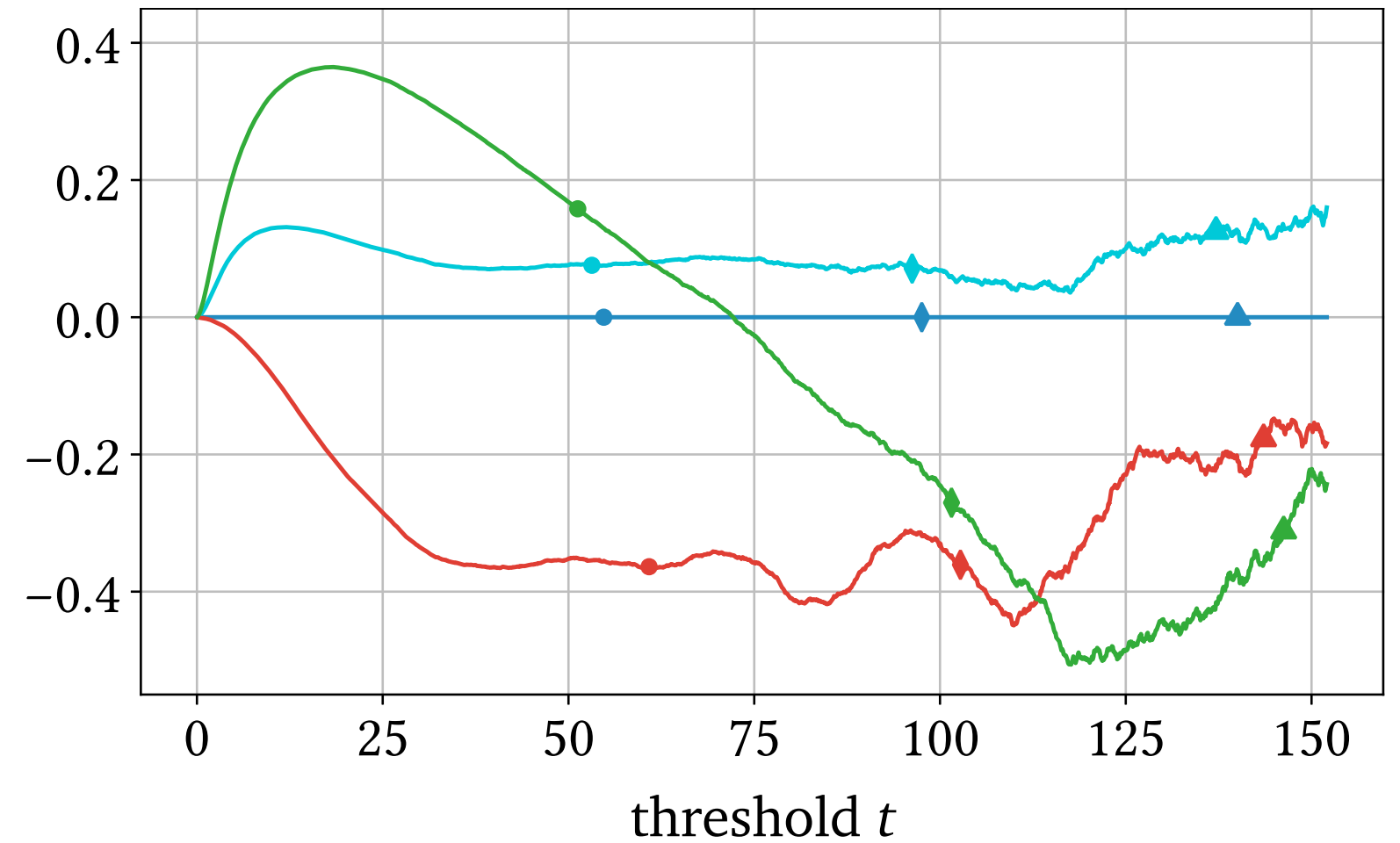


Boost large jobs?

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95

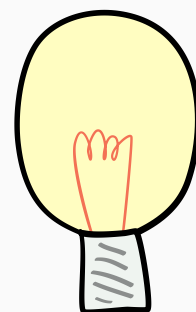


- 90th percentile

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

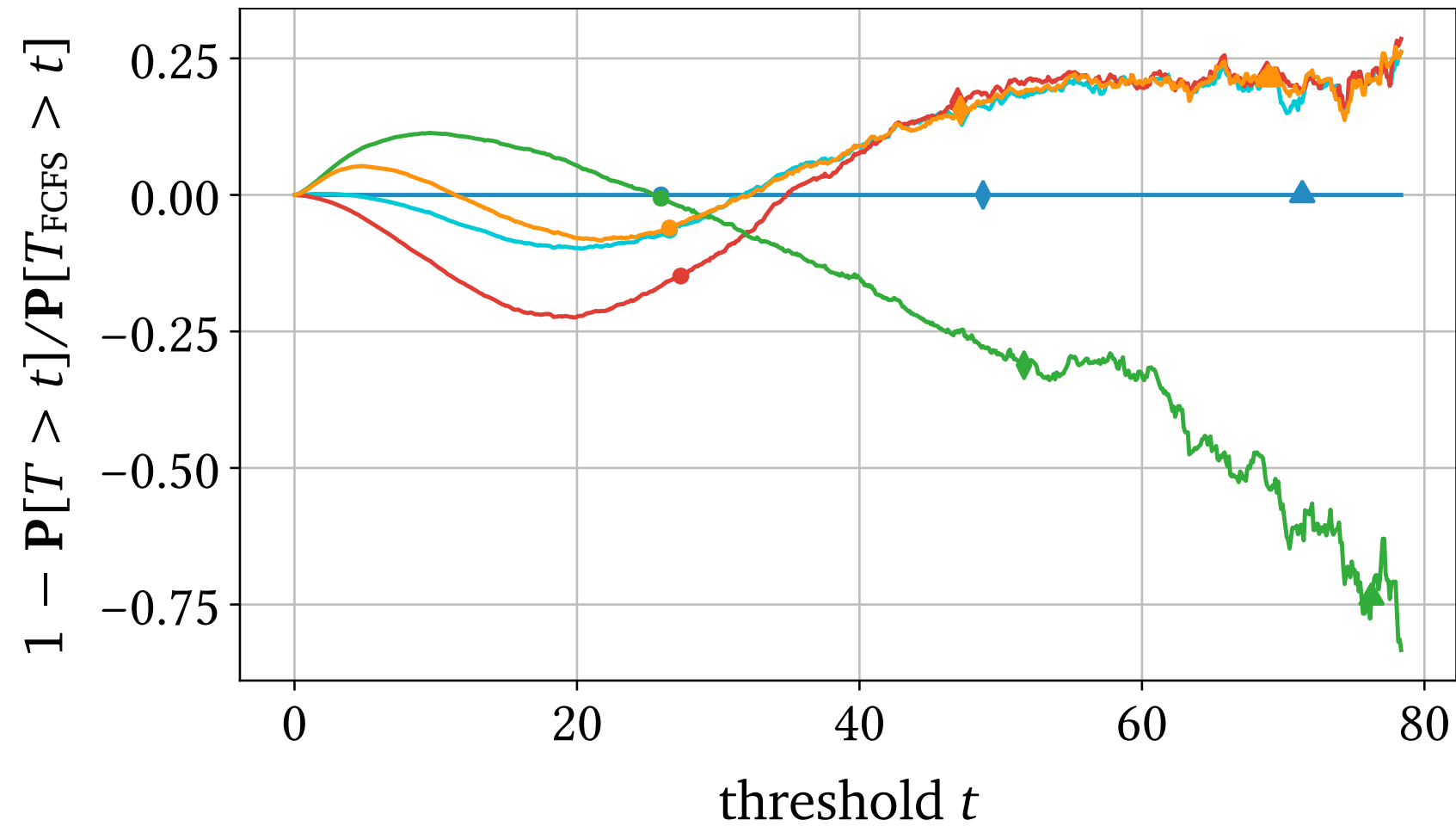
when $b(s) = (k-1)s$ 90th p

compare to
bin packing

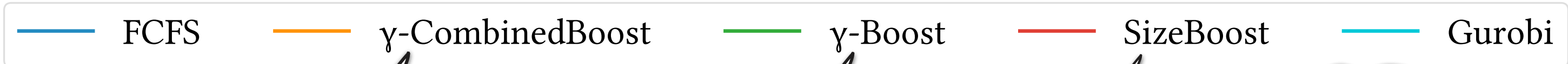
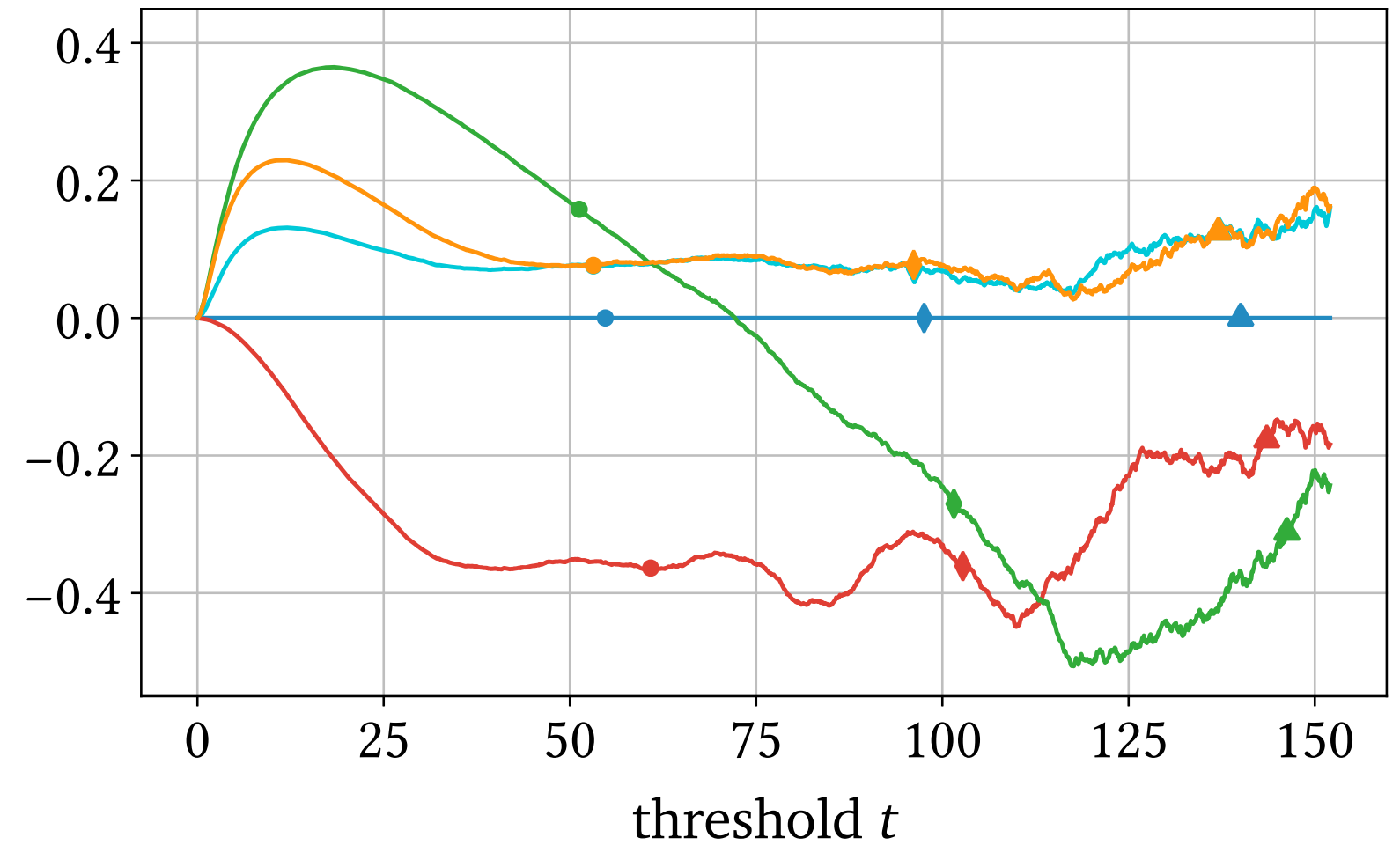


Boost large jobs?

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



99th percentile

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}} + (k-1)s$$

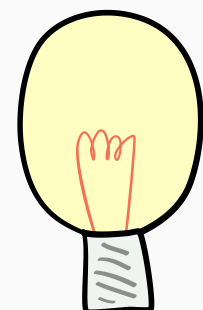
mean

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

9.9th percentile

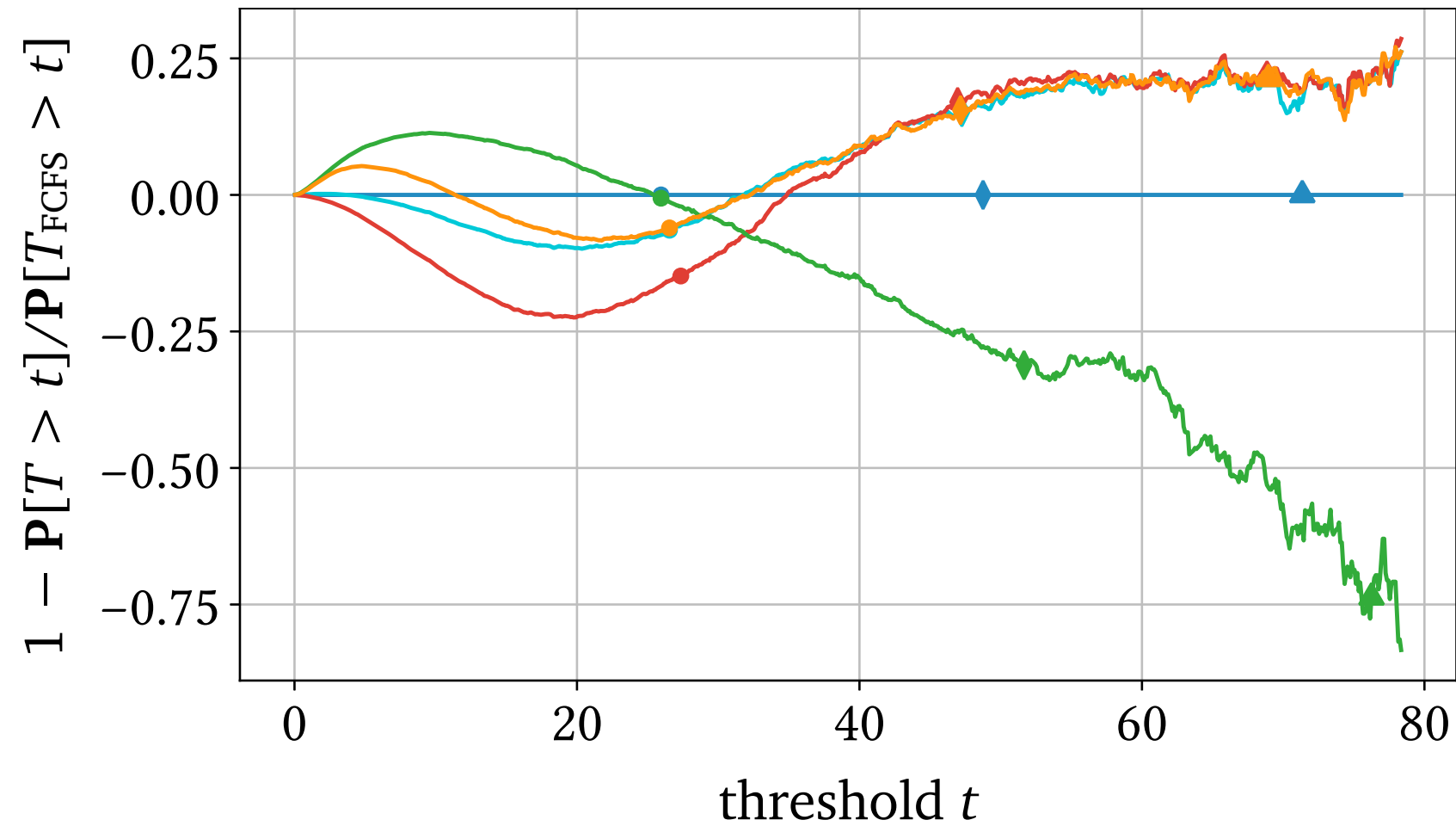
$$b(s) = (k-1)s$$

compare to bin packing

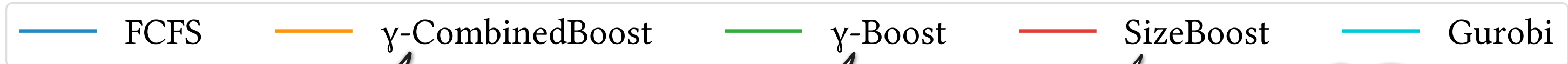
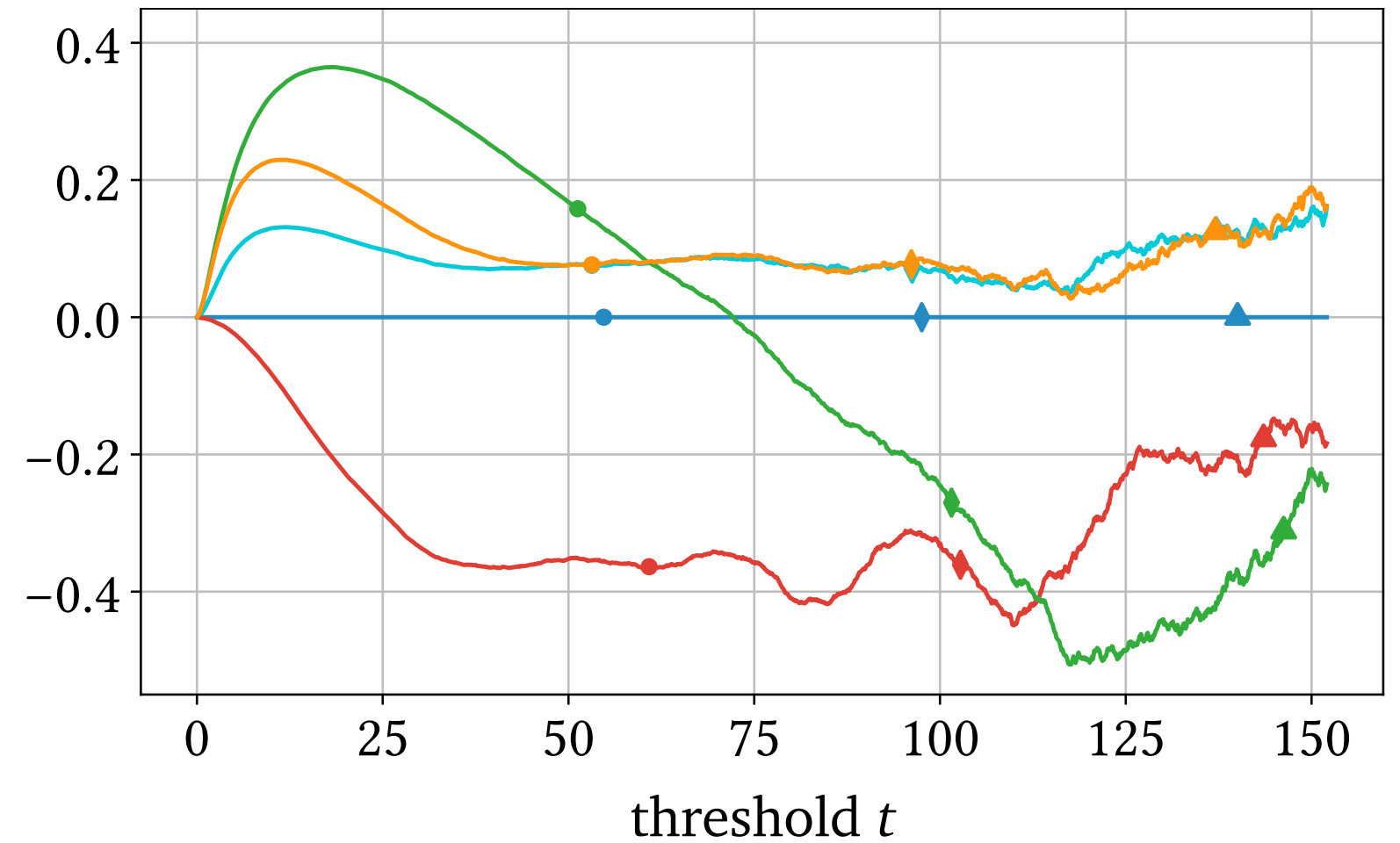


Boost large jobs?

$k = 10$ servers, load 0.8



$k = 10$ servers, load 0.95



99th percentile

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}} + (k-1)s$$

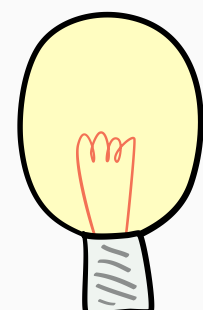
99th percentile

$$b(s) = \frac{1}{\gamma} \log \frac{1}{1 - e^{-\gamma s}}$$

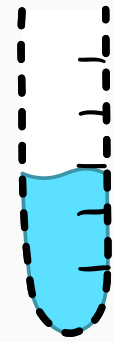
99th percentile

$$b(s) = (k-1)s$$

compare to bin packing

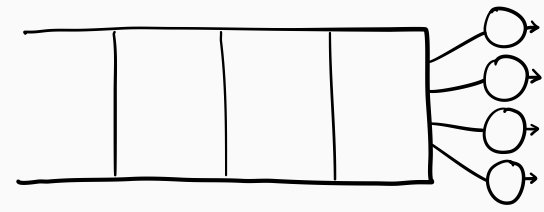


Boost large *and* small jobs!?



Part I

Handling job size uncertainty



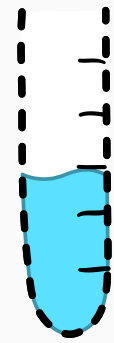
Part II

Analyzing multiserver scheduling



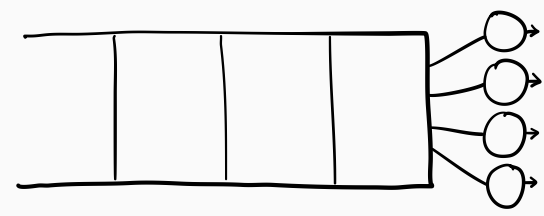
Part III

Optimizing tail metrics



Part I

Handling job size uncertainty



Part II

Analyzing multiserver scheduling

TCS for **Queueing**

*Approximation algorithms for
smoothed tail metric $\mathbb{E}[e^{\gamma T}]$?*



Part III

Optimizing tail metrics



Part I

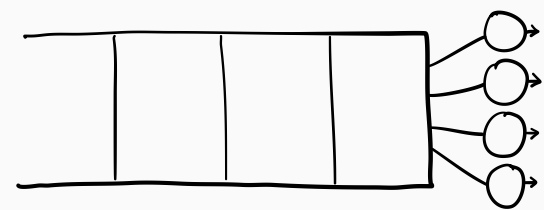
Handling job size uncertainty

Queueing for TCS

Use WINE to analyze Gittins with arbitrary release dates?

Queueing for TCS

Use WINE to analyze SRPT-*k* with arbitrary release dates?

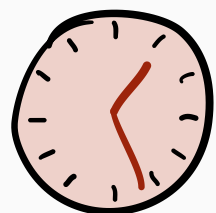


Part II

Analyzing multiserver scheduling

TCS for Queueing

Approximation algorithms for smoothed tail metric $\mathbf{E}[e^{\gamma T}]$?



Part III

Optimizing tail metrics