# Optimizing Tail Latency in Queues:

Generating Functions, the Gittins Index, and Gurobi

Ziv Scully

Cornell ORIE

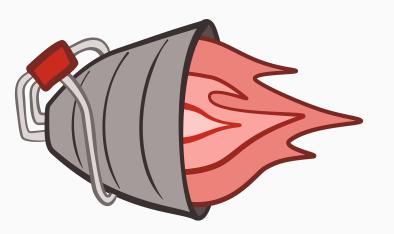
Joint work with

George Yu

Cornell ORIE

Amit Harley

Cornell CAM



# Optimizing Tail Latency in Queues:

Generating Functions, the Gittins Index, and Gurobi

SIGMETRICS 2024

Ziv Scully

Cornell ORIE

Joint work with

George Yu

Cornell ORIE

Amit Harley

Cornell CAM



## Optimizing Tail Latency in Queues:

Generating Functions, the Gittins Index, and Gurobi

SIGMETRICS 2024

SIGMETRICS 2025

Ziv Scully

Cornell ORIE

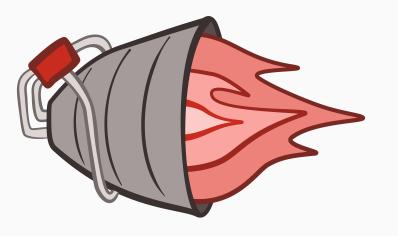
Joint work with

George Yu

Cornell ORIE

Amit Harley

Cornell CAM



# Optimizing Tail Latency in Queues:

Generating Functions, the Gittins Index, and Gurobi

SIGMETRICS 2024)

SIGMETRICS 2025

SIGMETRICS 2026

Ziv Scully

Cornell ORIE

Joint work with

George Yu Amit Harlev Cornell ORIE

Cornell CAM



## Optimizing Tail Latency in Queues:

Generating Functions, the Gittins Index, and Gurobi

SIGMETRICS 2024

SIGMETRICS 2025

SIGMETRICS 2026

Ziv Scully

SIGMETRICS 2024 best paper, 2024 INFORMS APS student paper competition finalist

Joint work wit

George Yu

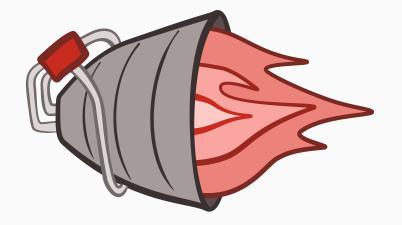
Cornell ORIE

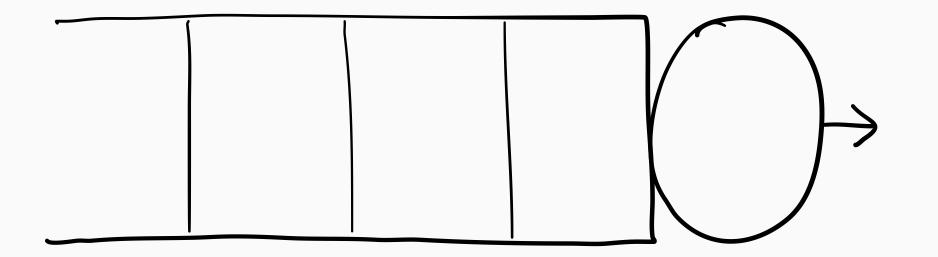
Amit Harlev

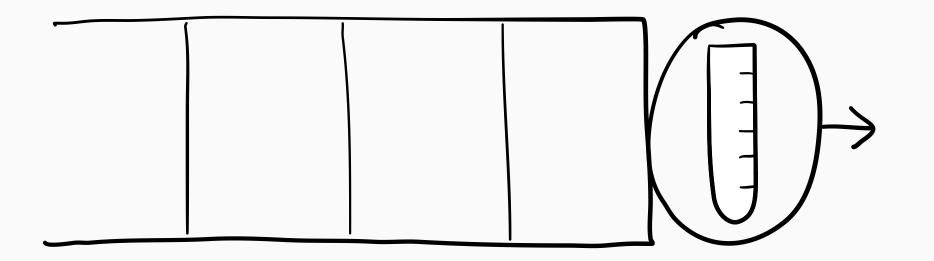
Cornell CAM

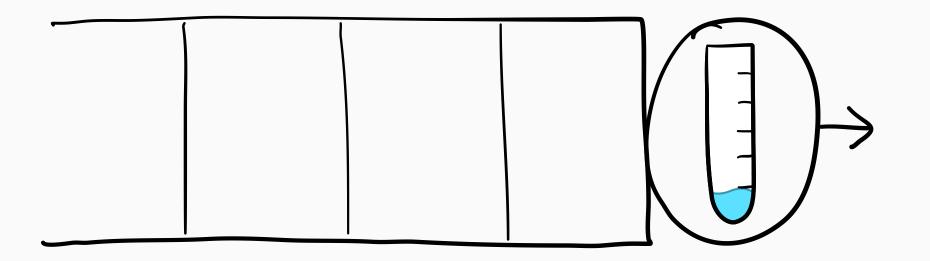
Reevu Adakroy

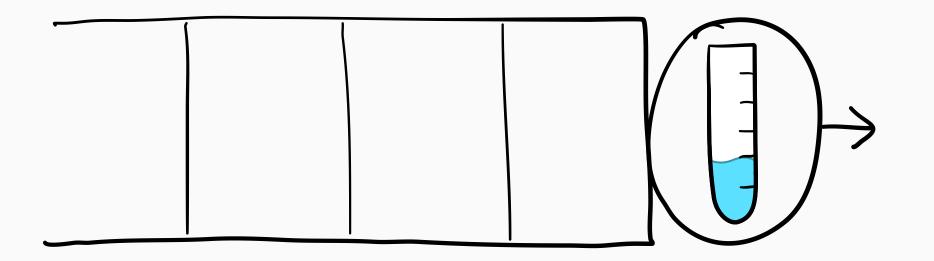
Cornell ORIE

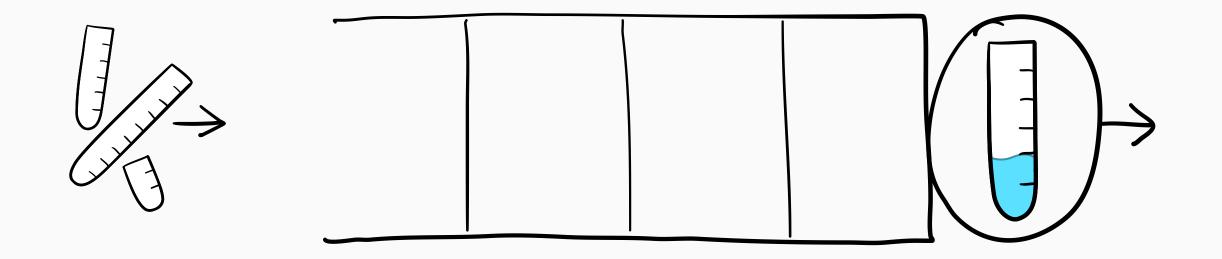


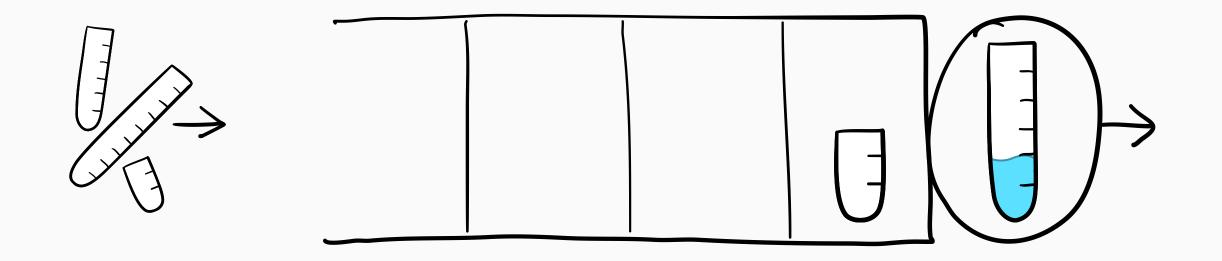


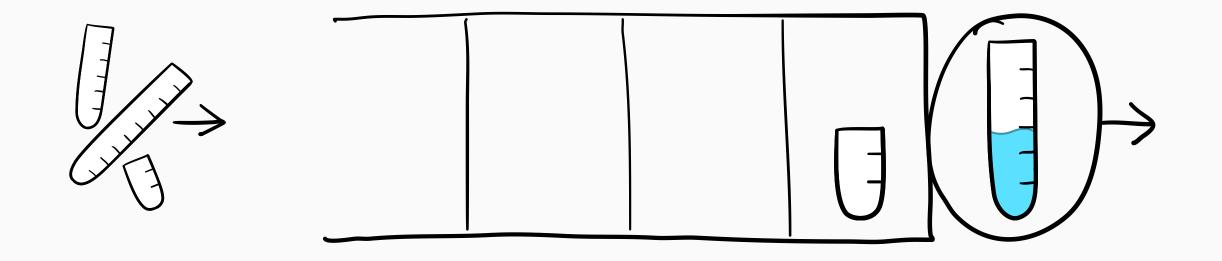


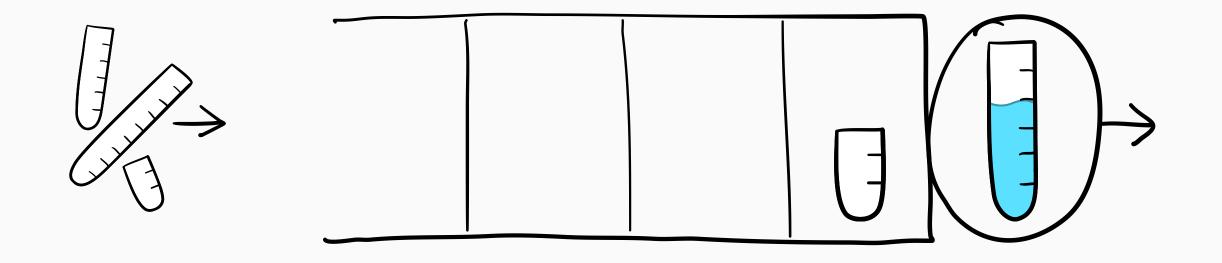


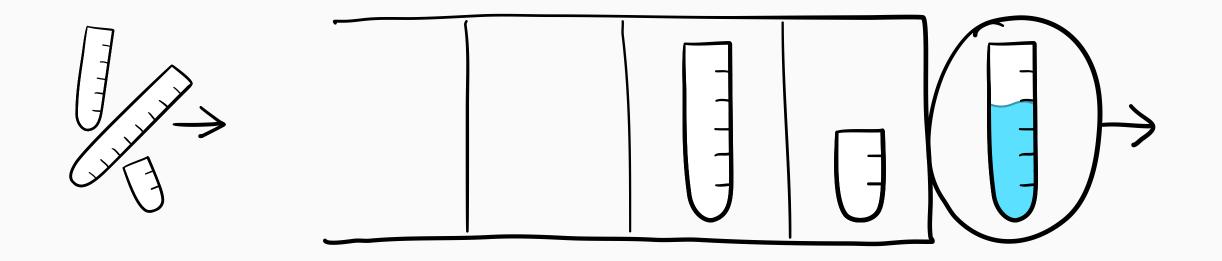


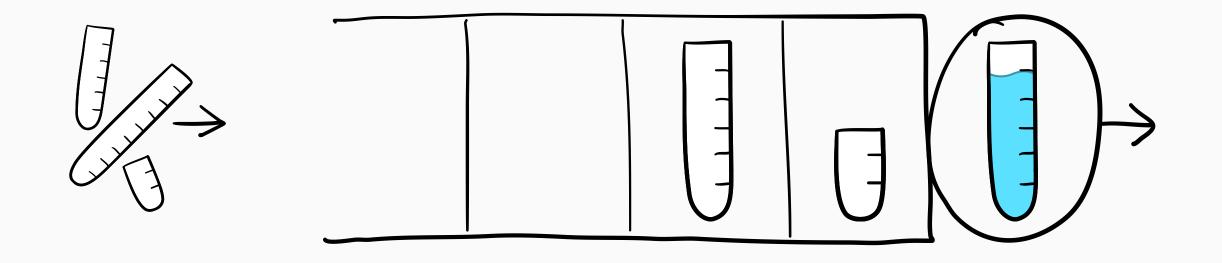


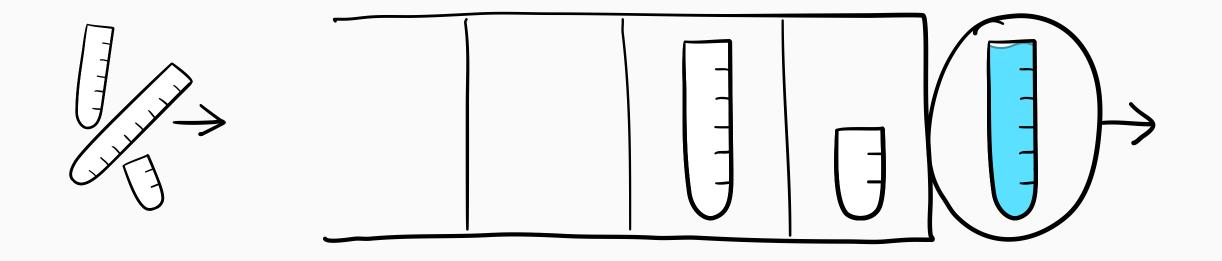


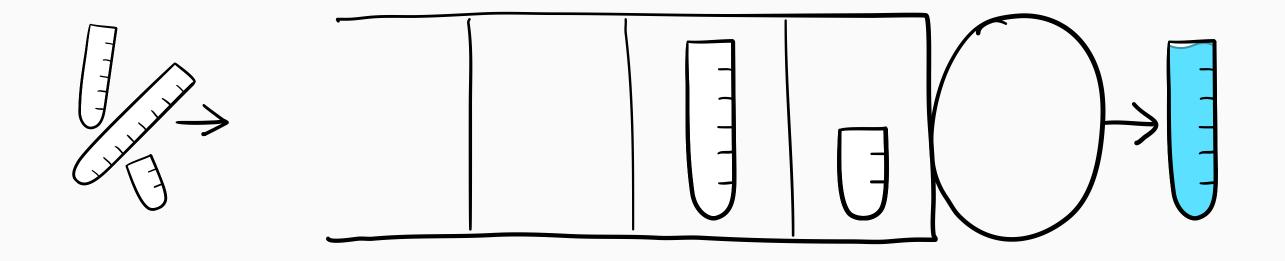


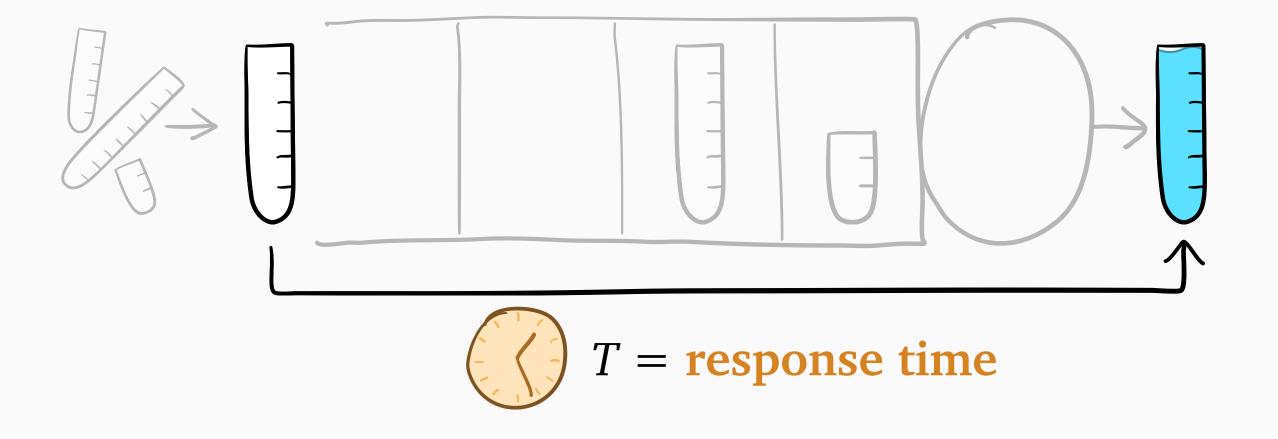


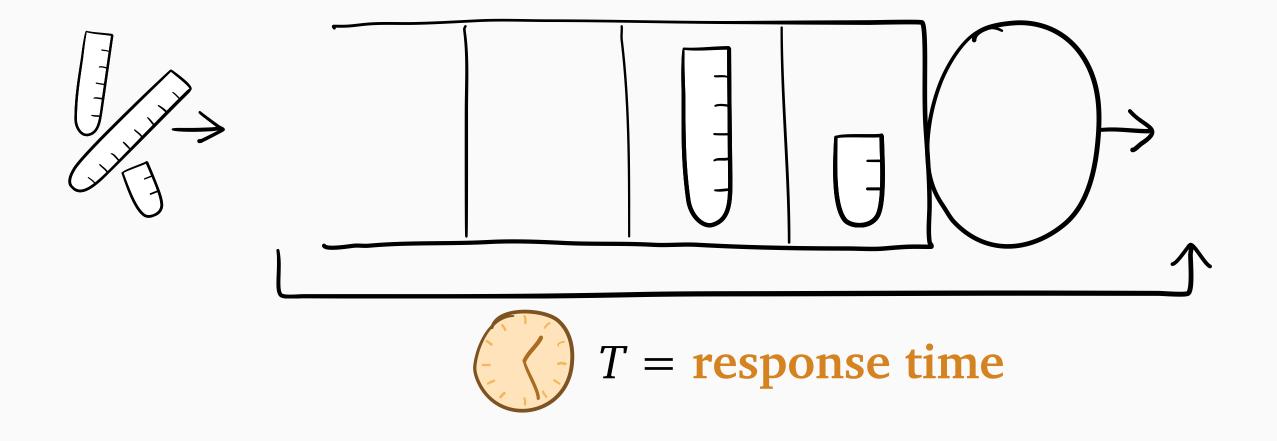


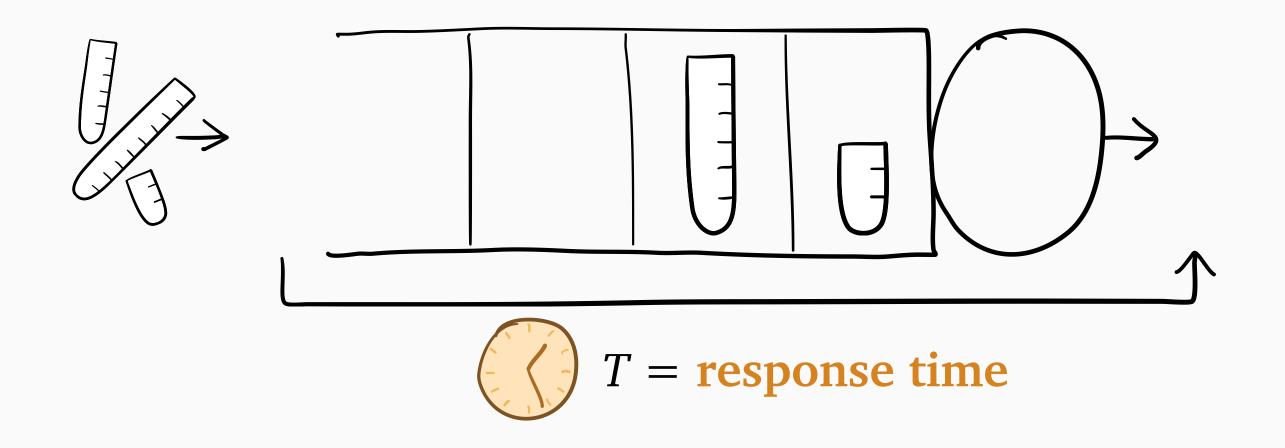




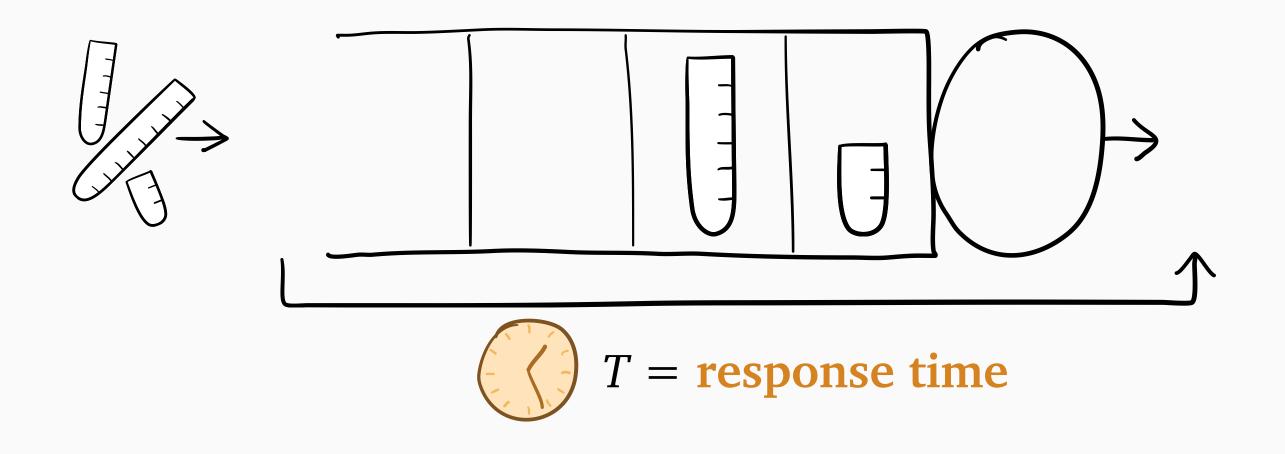


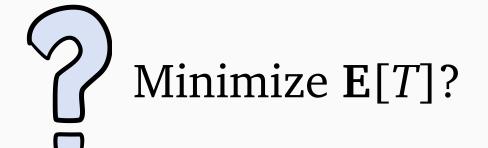


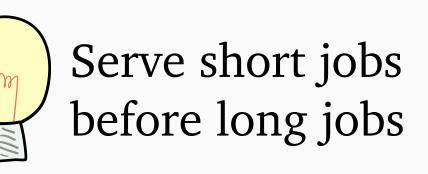


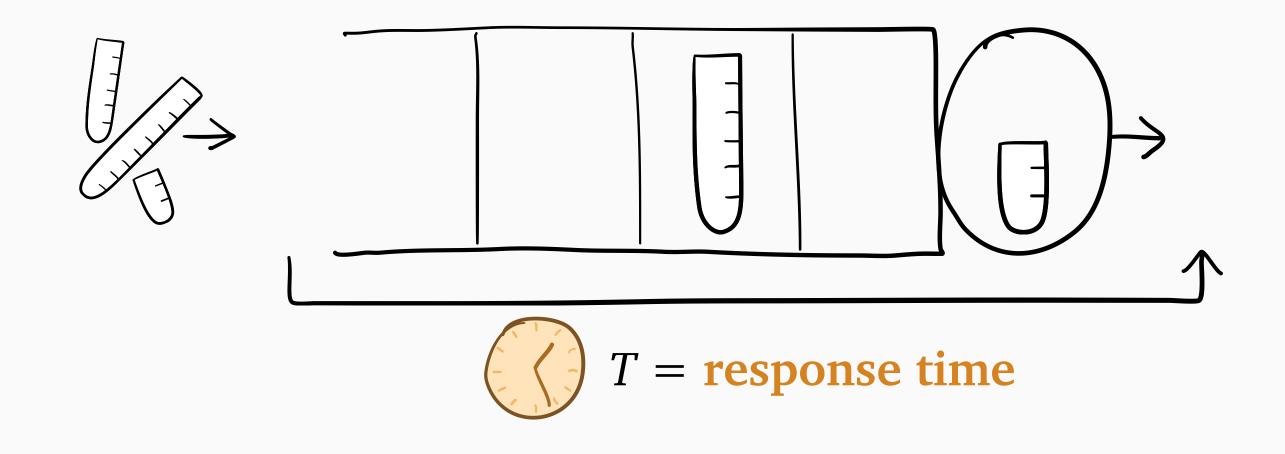


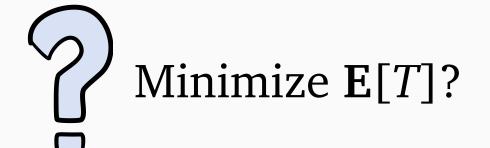


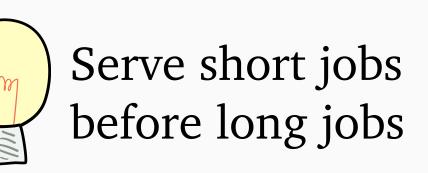


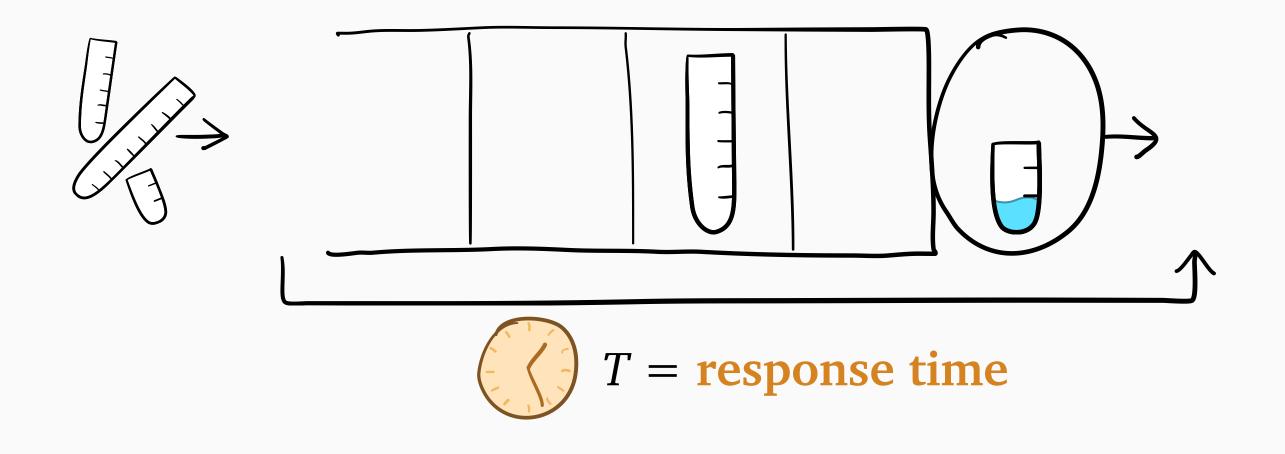


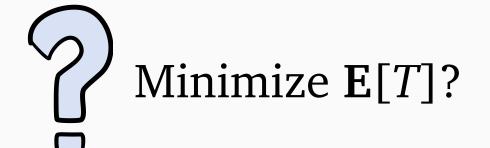


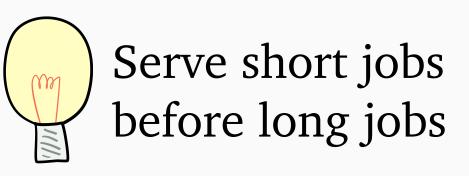


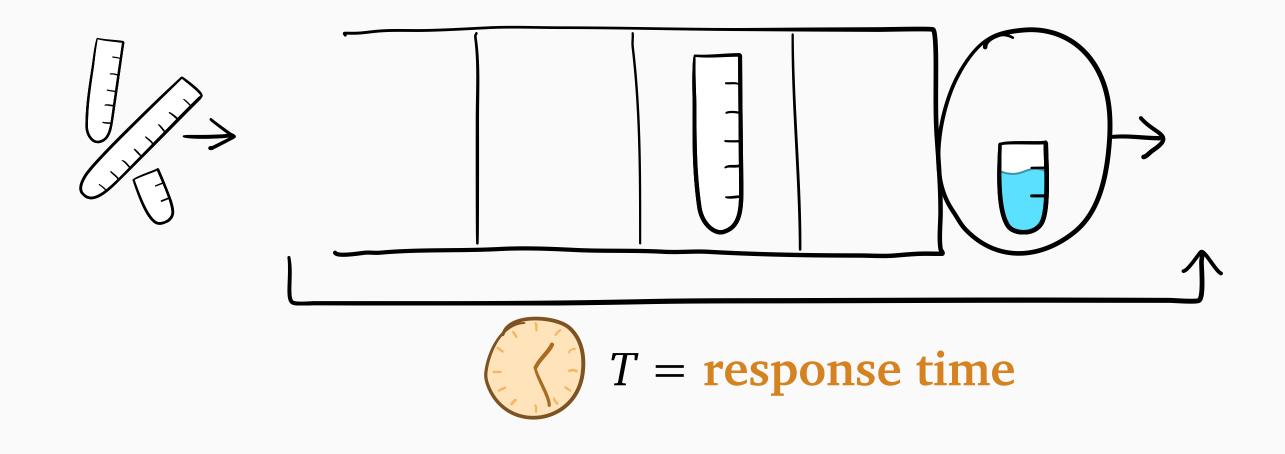


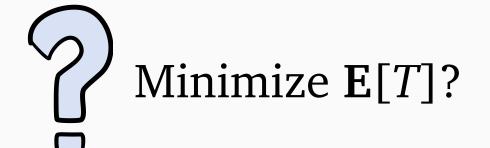


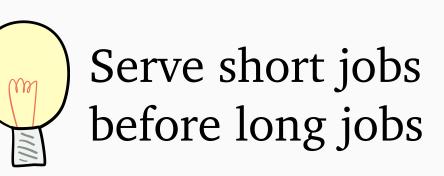


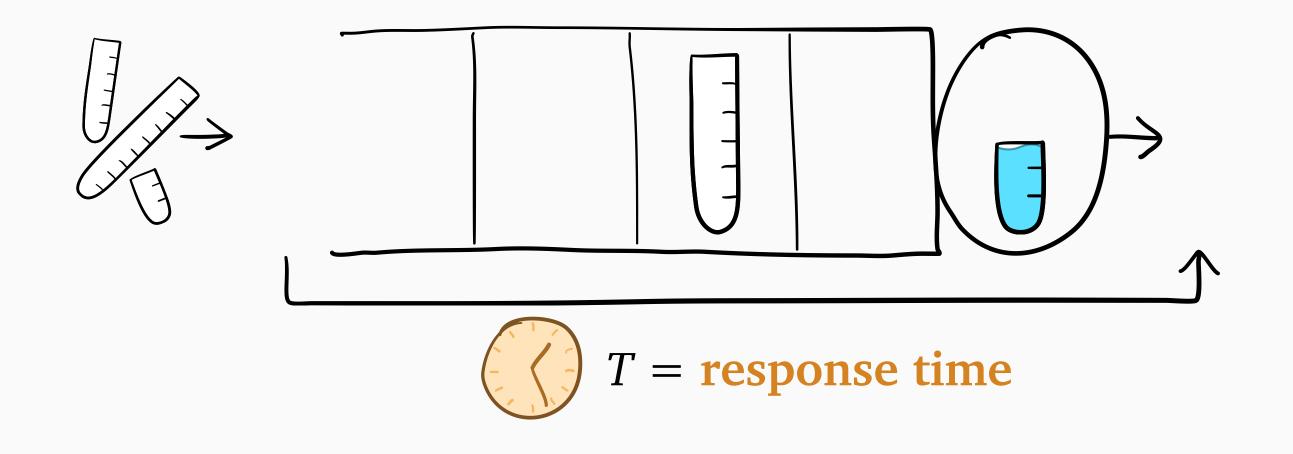


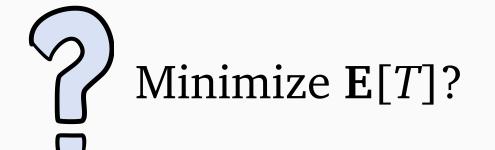


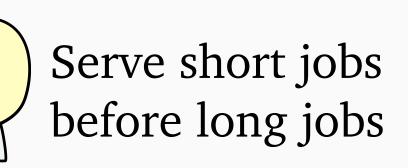


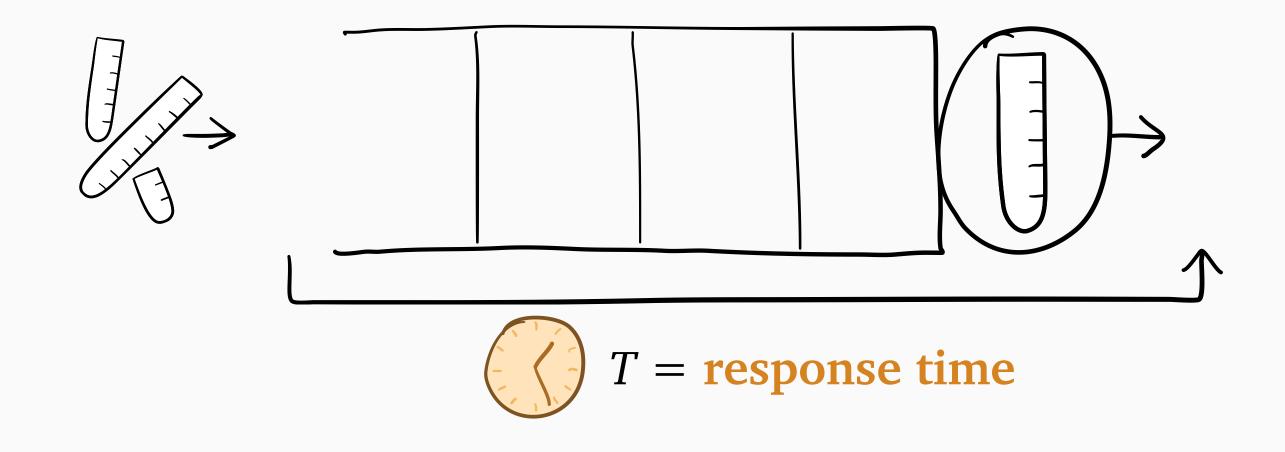


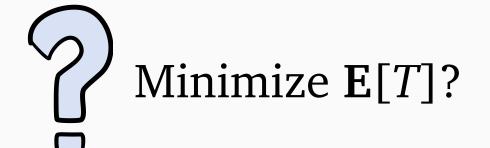


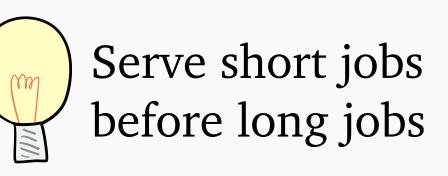


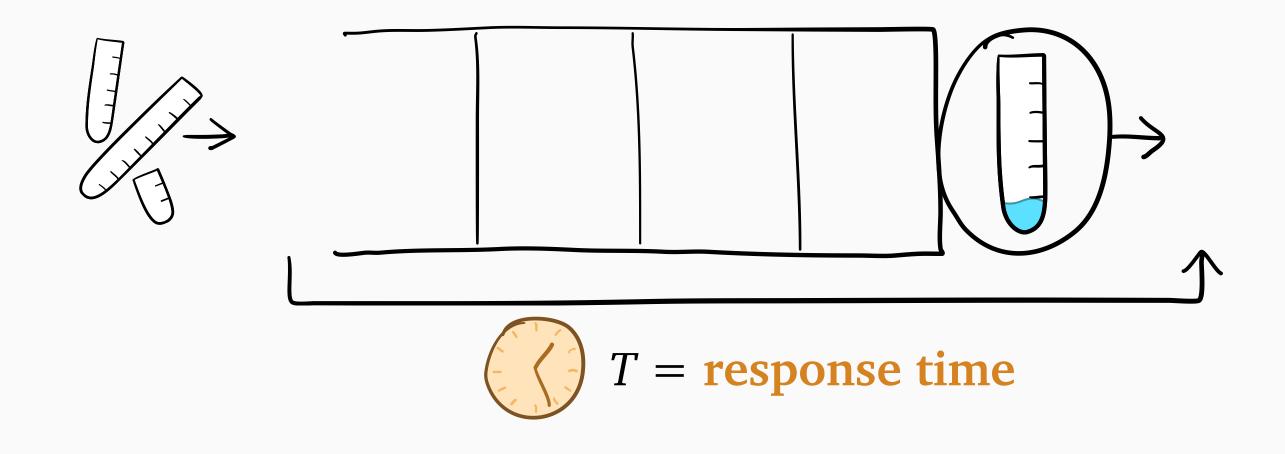


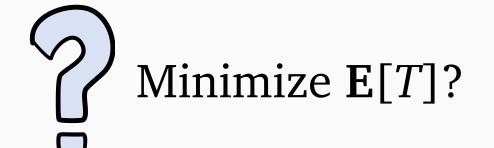


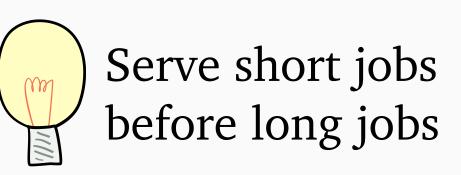


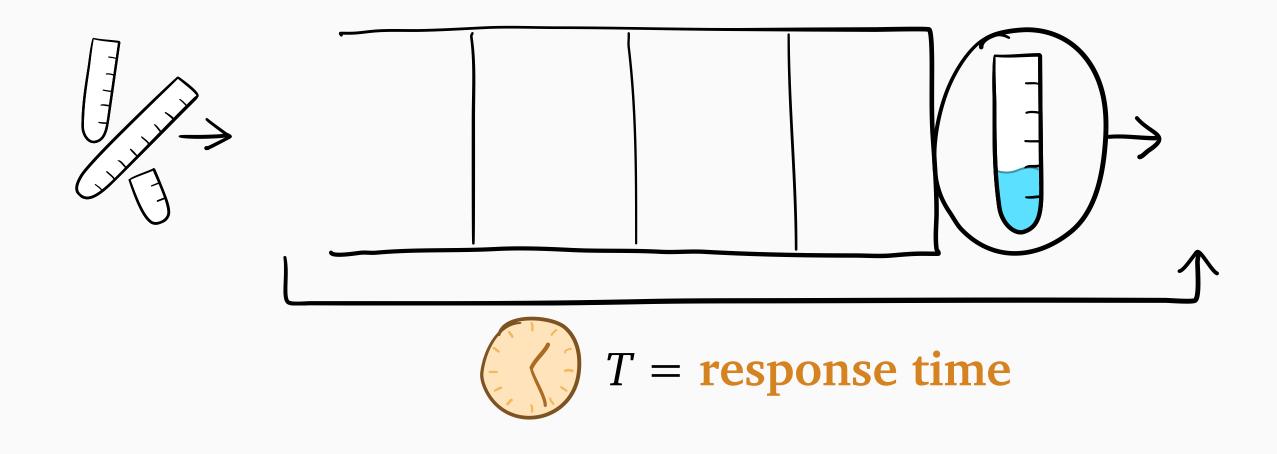


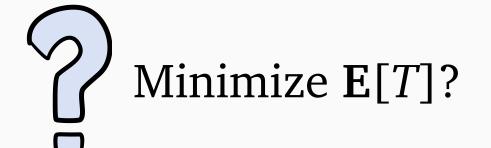


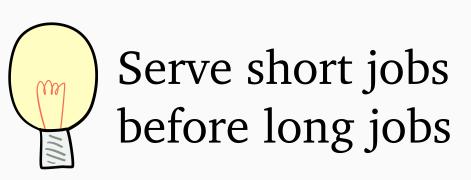


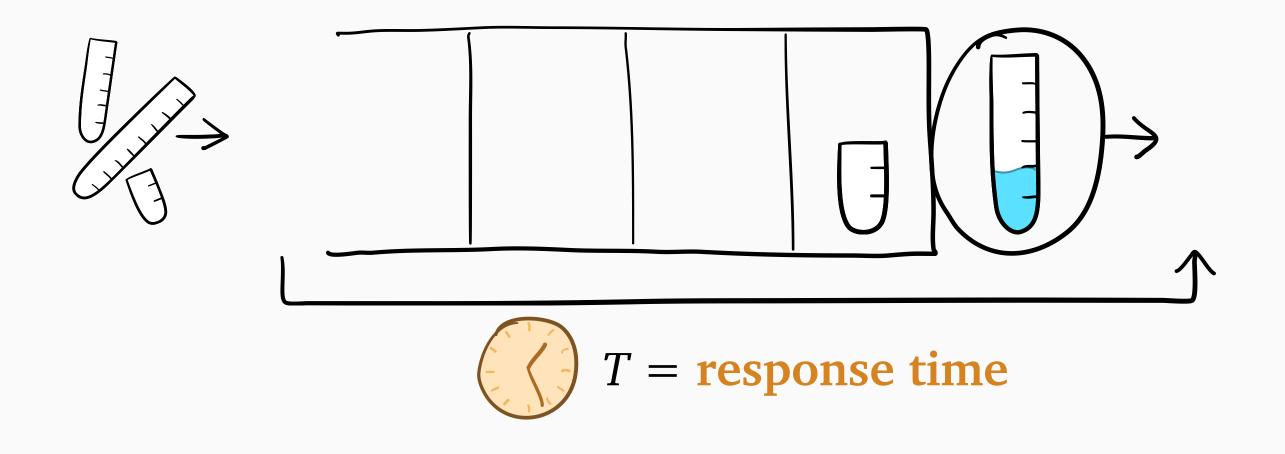


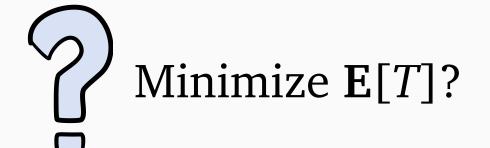


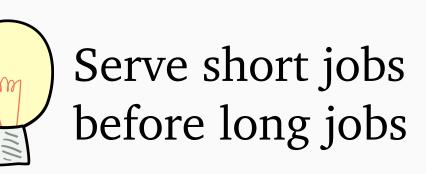


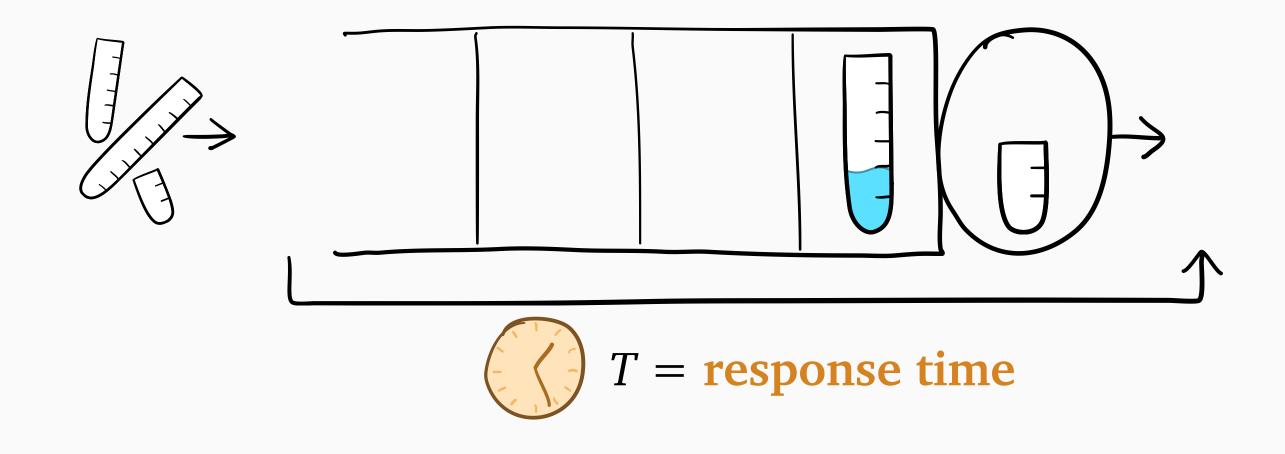


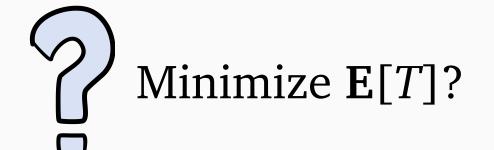


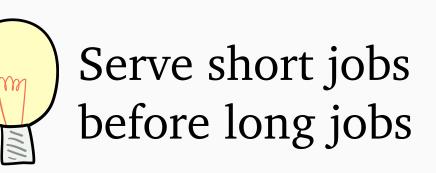


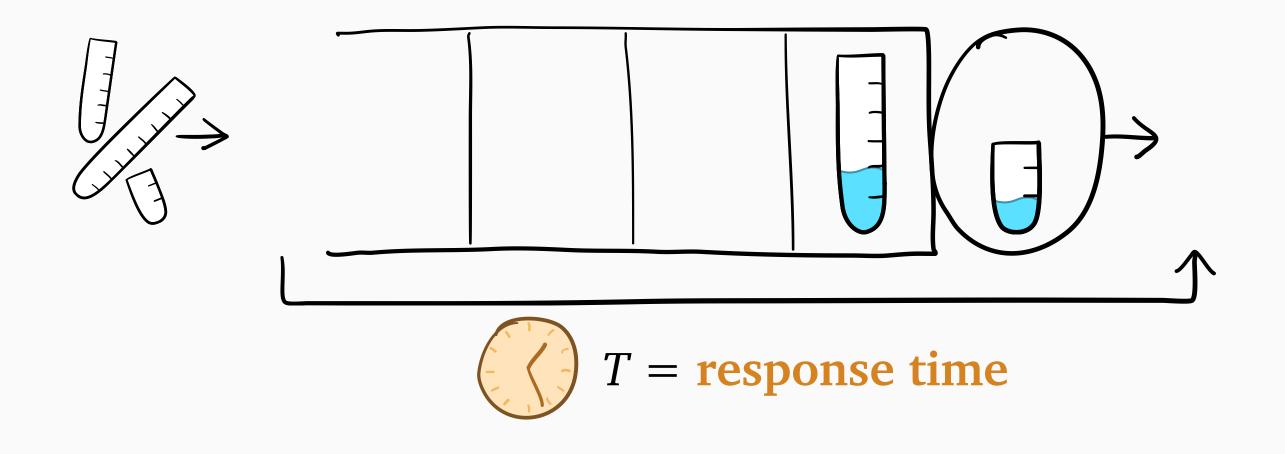


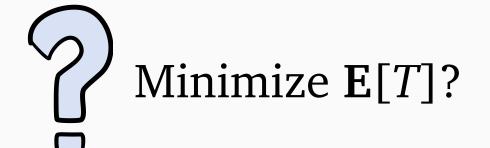


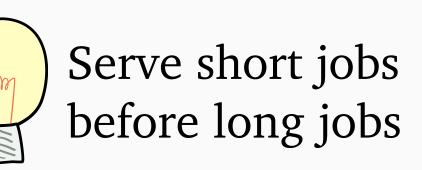


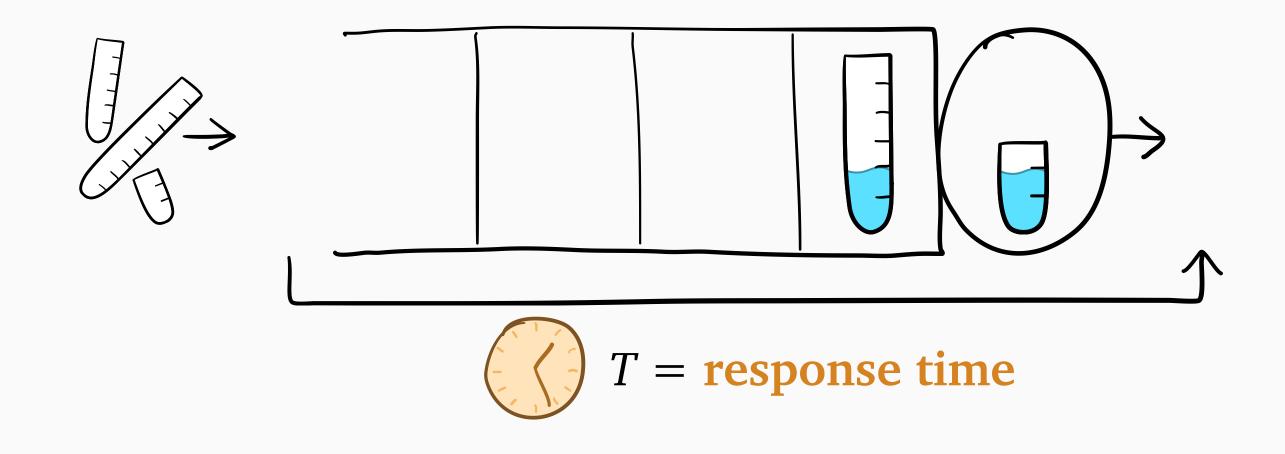


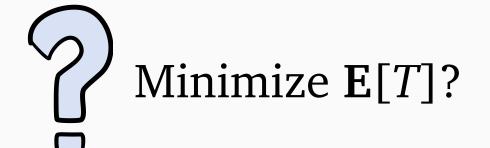




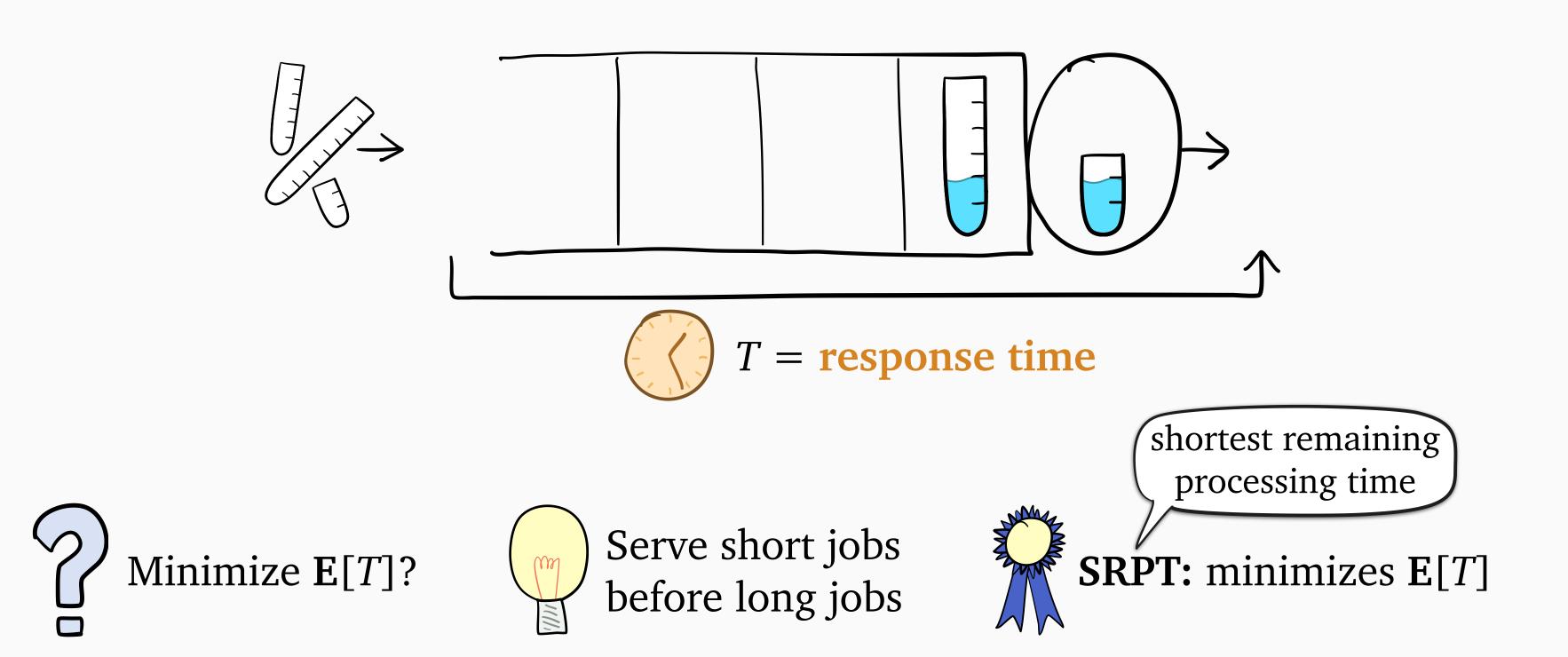




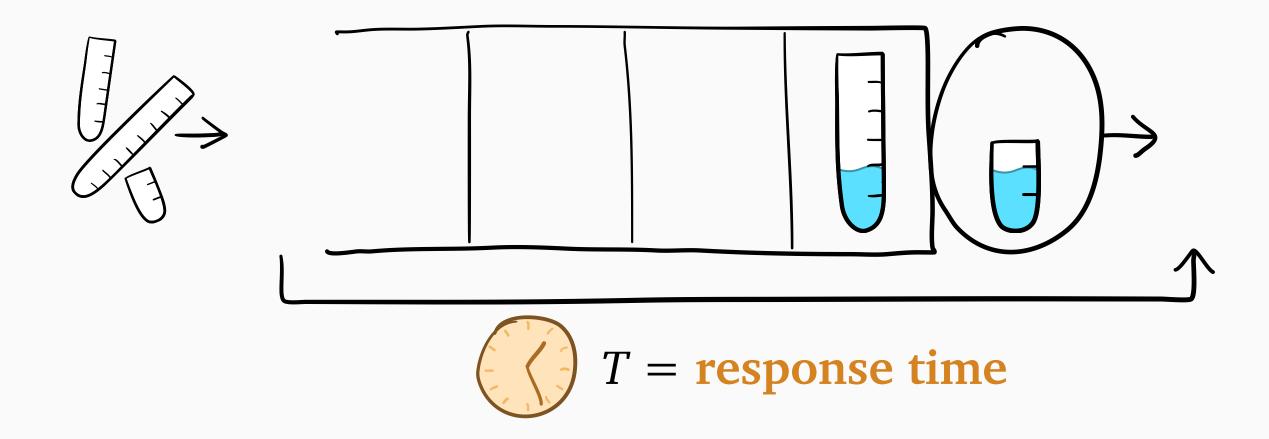






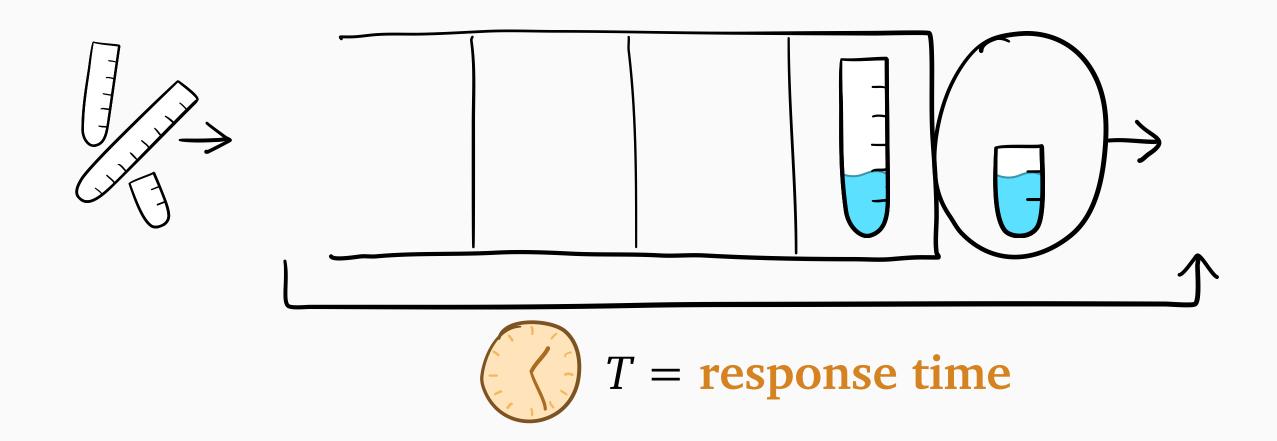


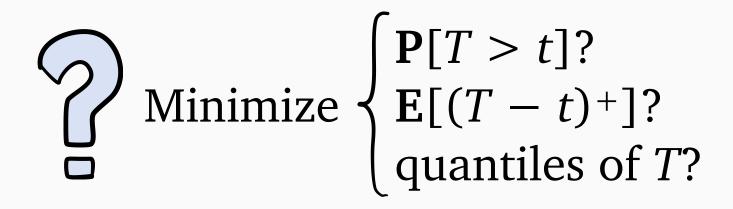
### Beyond the mean: tail metrics



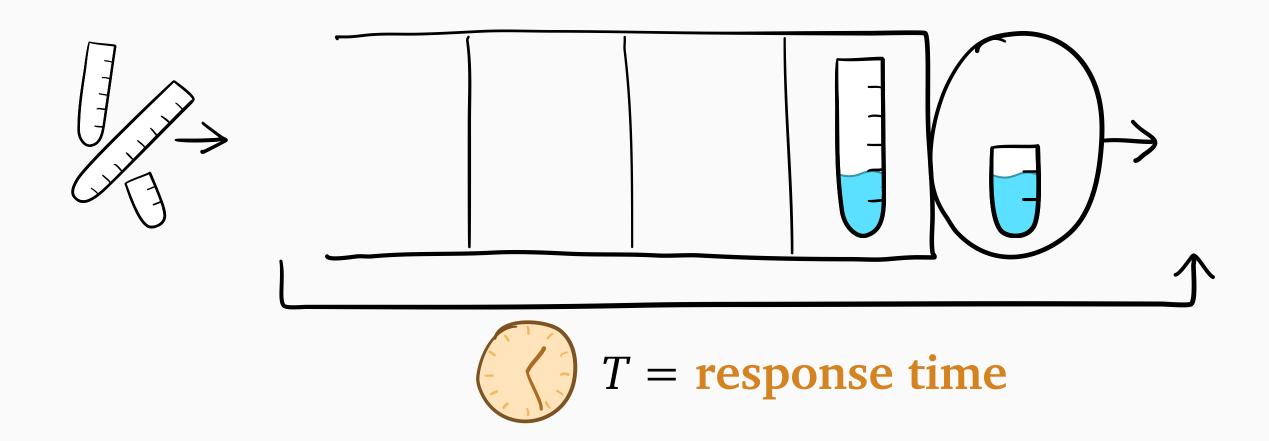
#### Beyond the mean: tail metrics

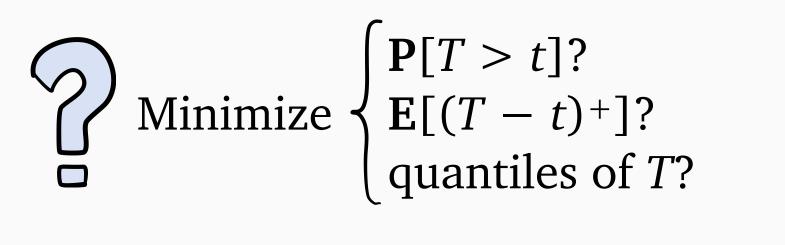
Minimize 
$$\begin{cases} \mathbf{P}[T > t]? \\ \mathbf{E}[(T - t)^{+}]? \\ \text{quantiles of } T? \end{cases}$$









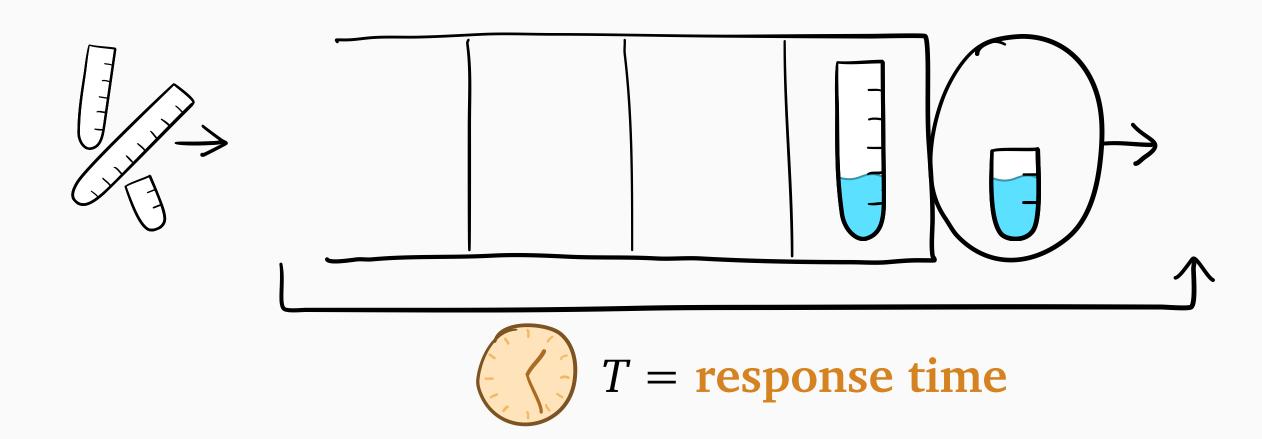


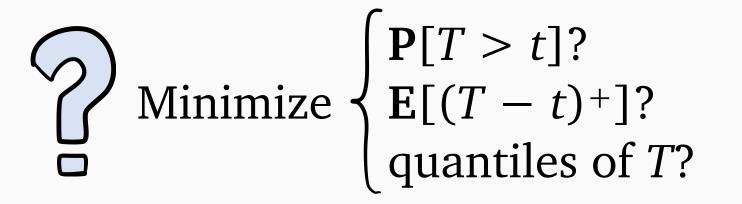


**Practice:** important



Theory: very hard



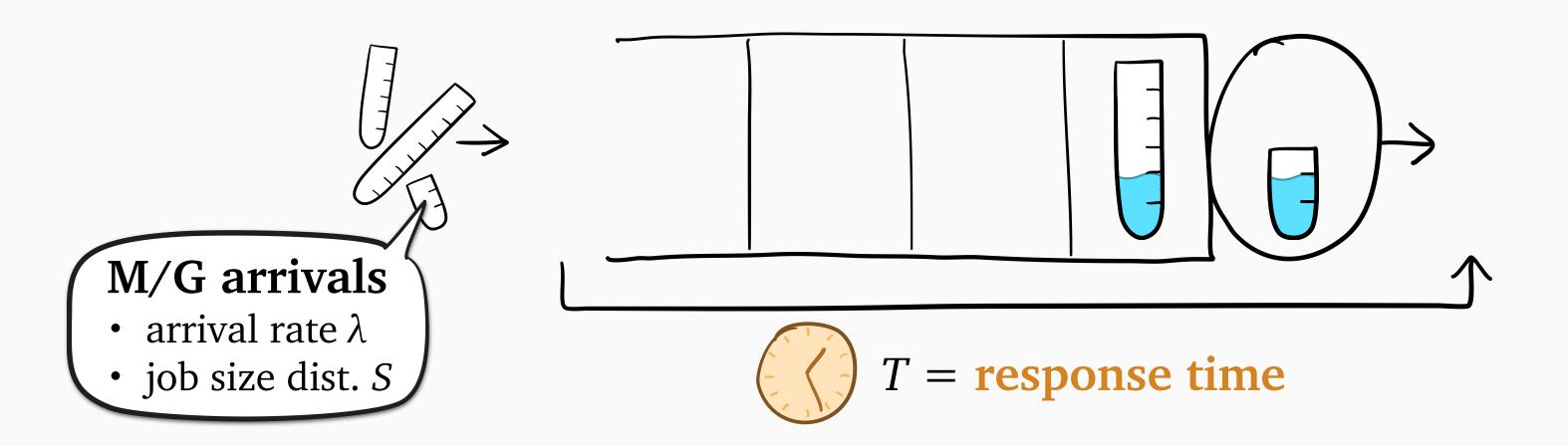


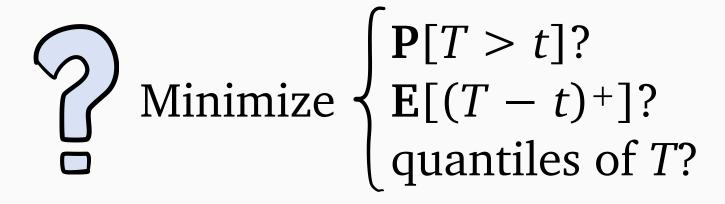


**Practice:** important



Theory: very hard







**Practice:** important



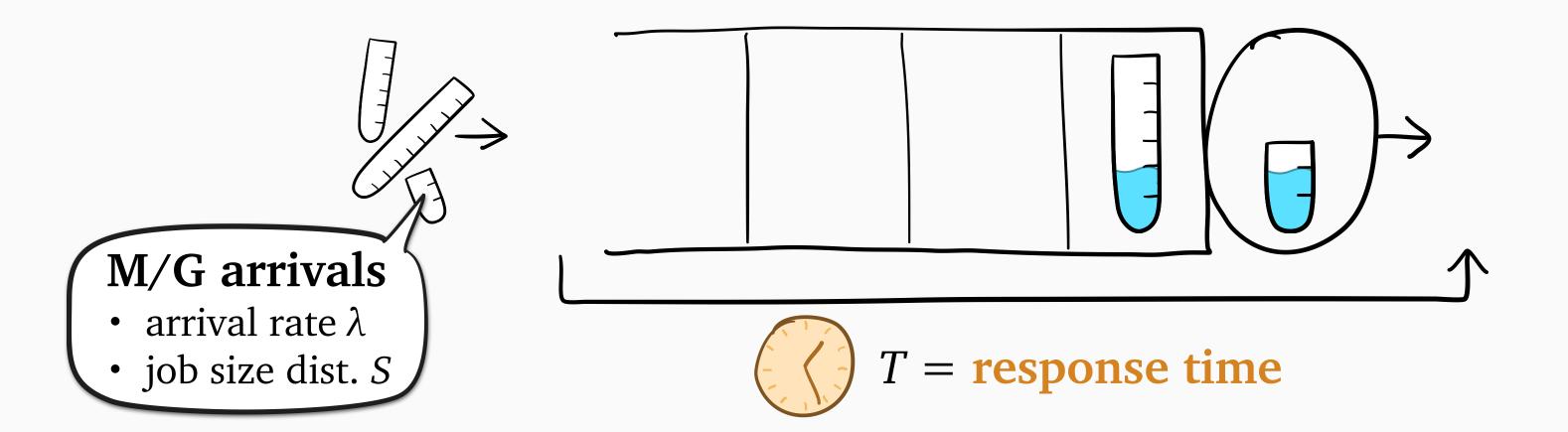
Theory: very hard

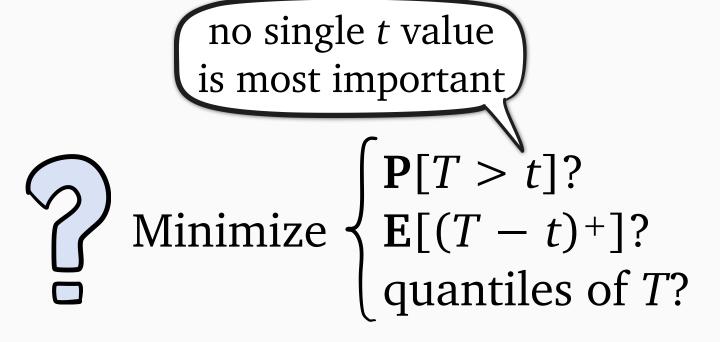


**Tractable:** 

study  $t \rightarrow \infty$ 

asymptotics







**Practice:** important



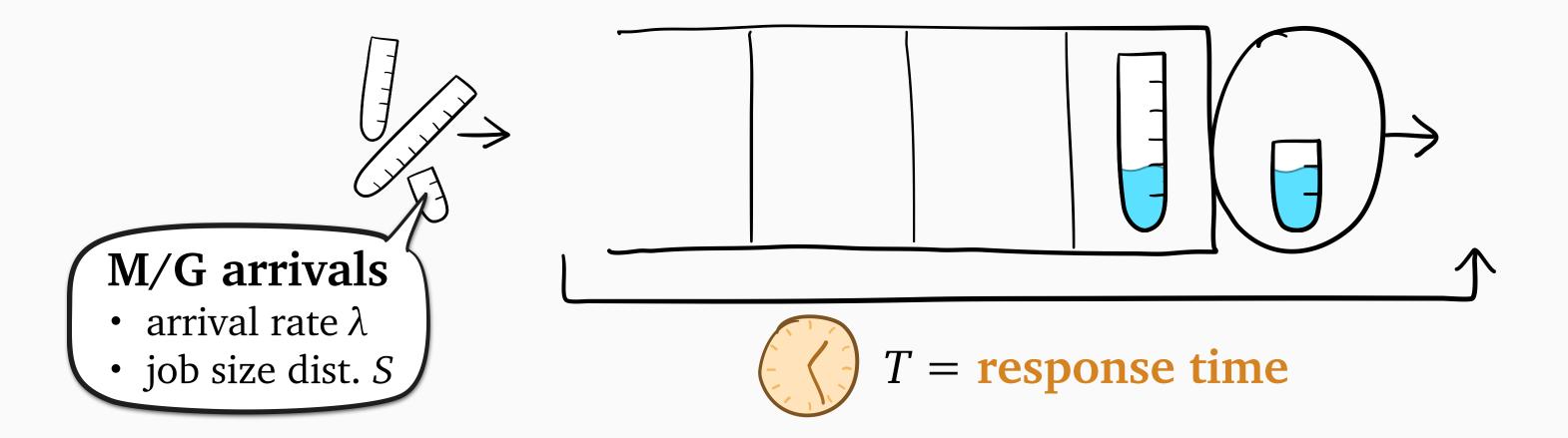
Theory: very hard

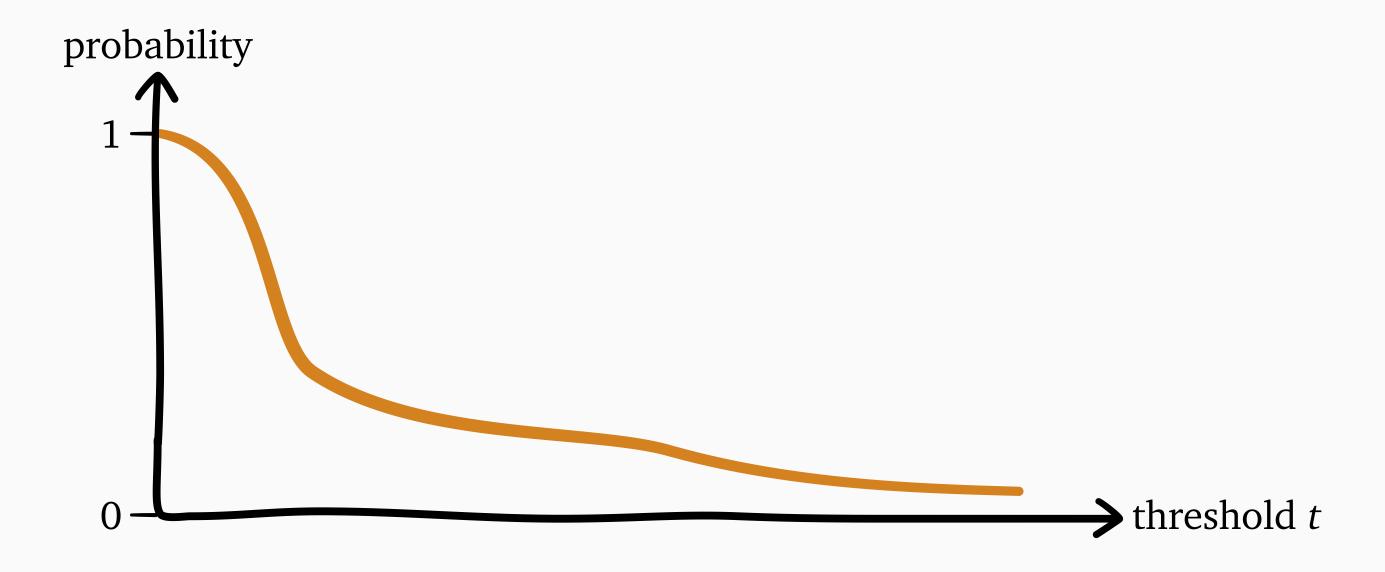


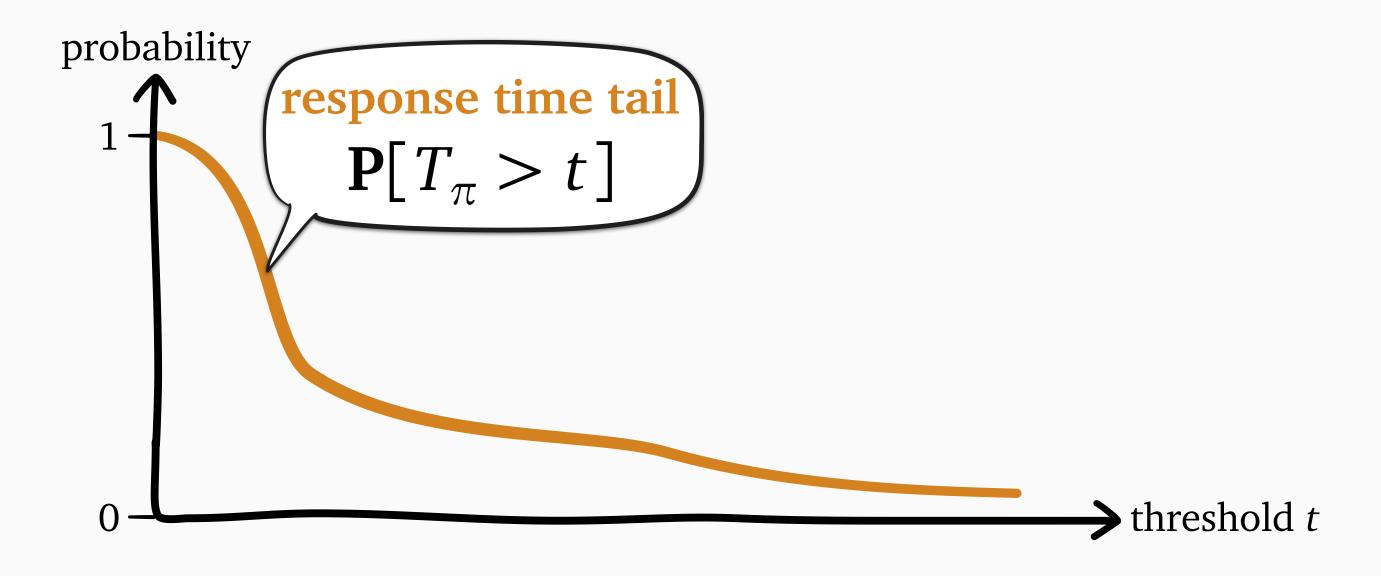
**Tractable:** 

study  $t \to \infty$ 

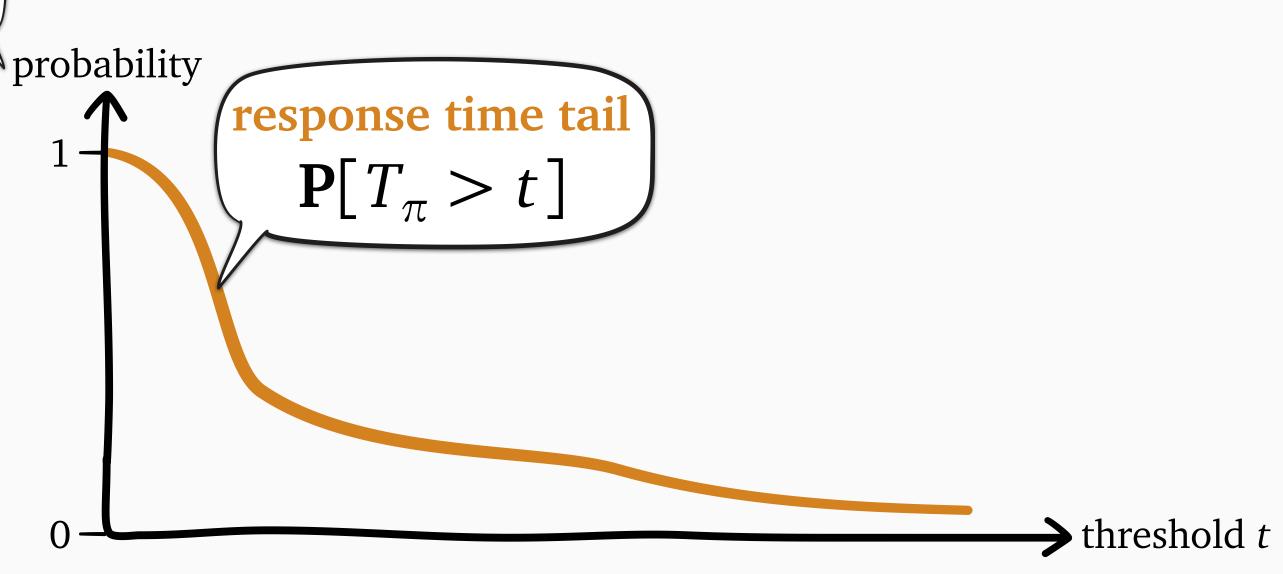
asymptotics



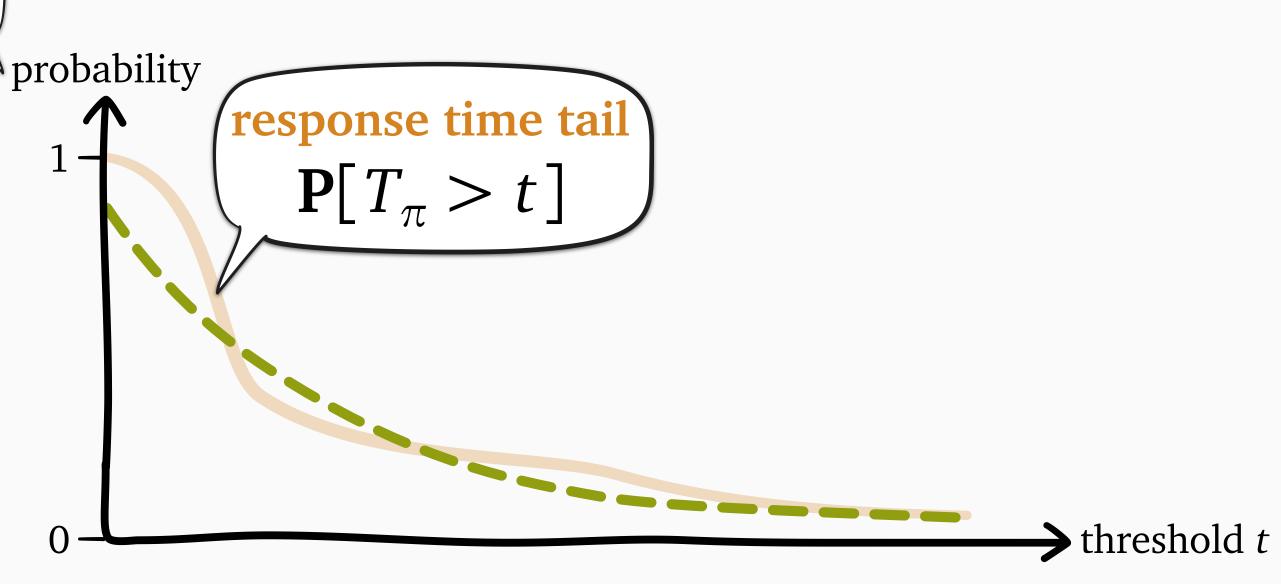




depends on policy  $\pi$ 



depends on policy  $\pi$ 



Asymptotic response time taily depends on when *S* is policy  $\pi$ light-tailed probability response time tail  $\mathbf{P}[T_{\pi} > t]$ asymptotic behavior  $C_{\pi}e^{-\gamma_{\pi}t}$ 

 $\rightarrow$  threshold t

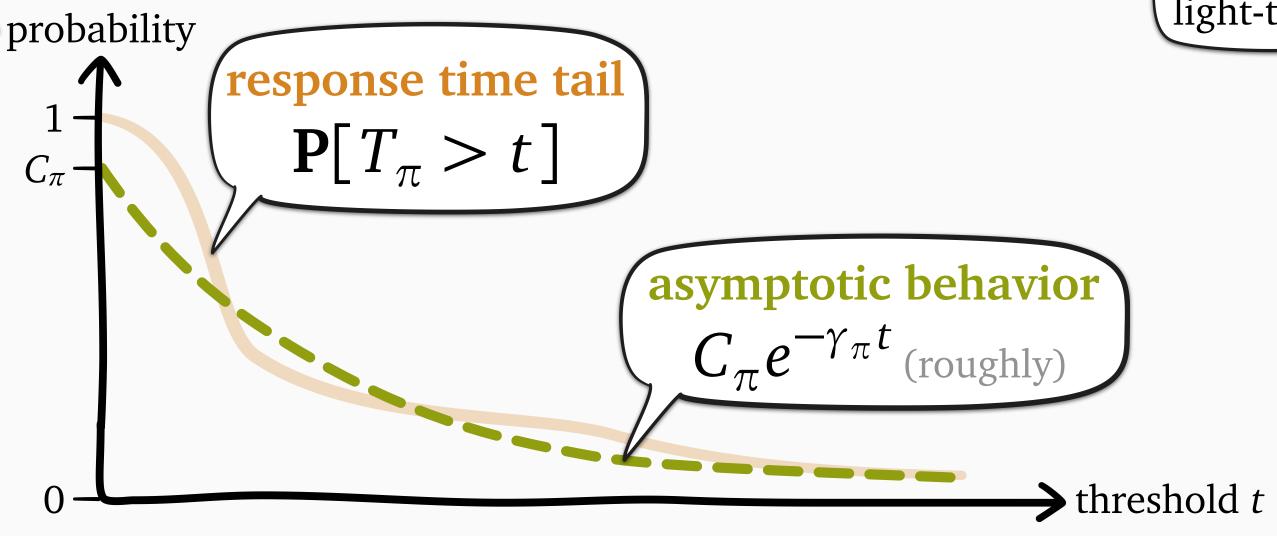
when *S* is light-tailed

 $\rightarrow$  threshold t

depends on policy  $\pi$ 

## Asymptotic response time taily

when *S* is light-tailed



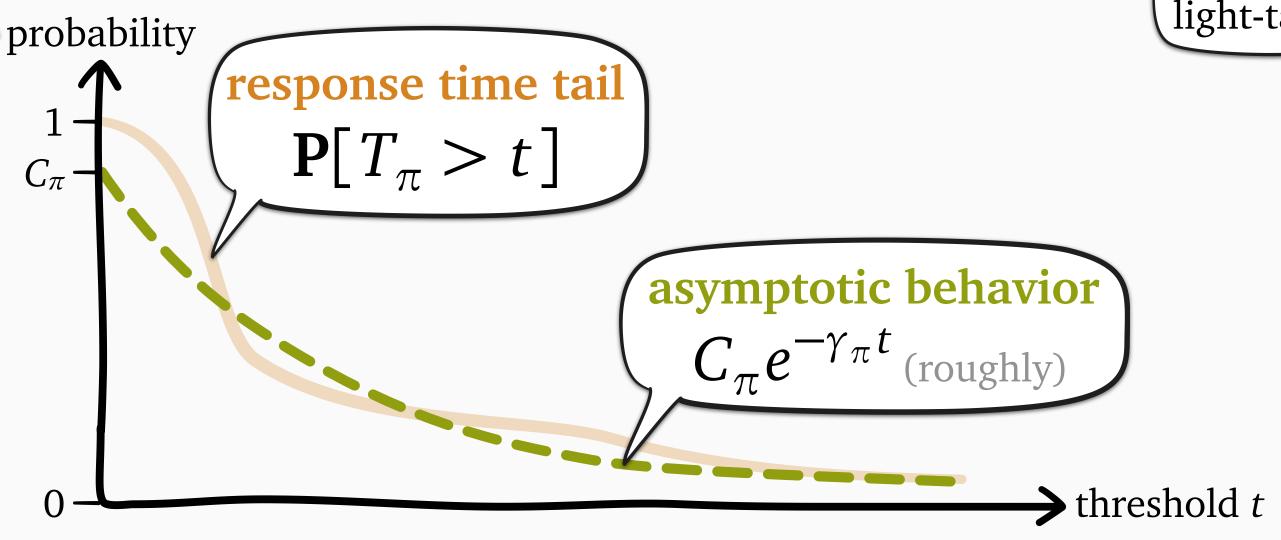
$$\gamma_{\pi} = decay \ rate \ of \ \pi$$

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

depends on policy  $\pi$ 

## Asymptotic response time taily

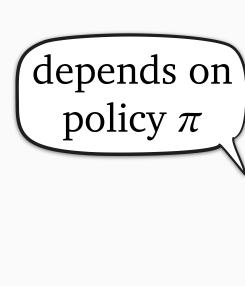
when *S* is light-tailed



Weak optimality: 
$$\leftarrow$$
 optimal  $\gamma_{\pi}$ 

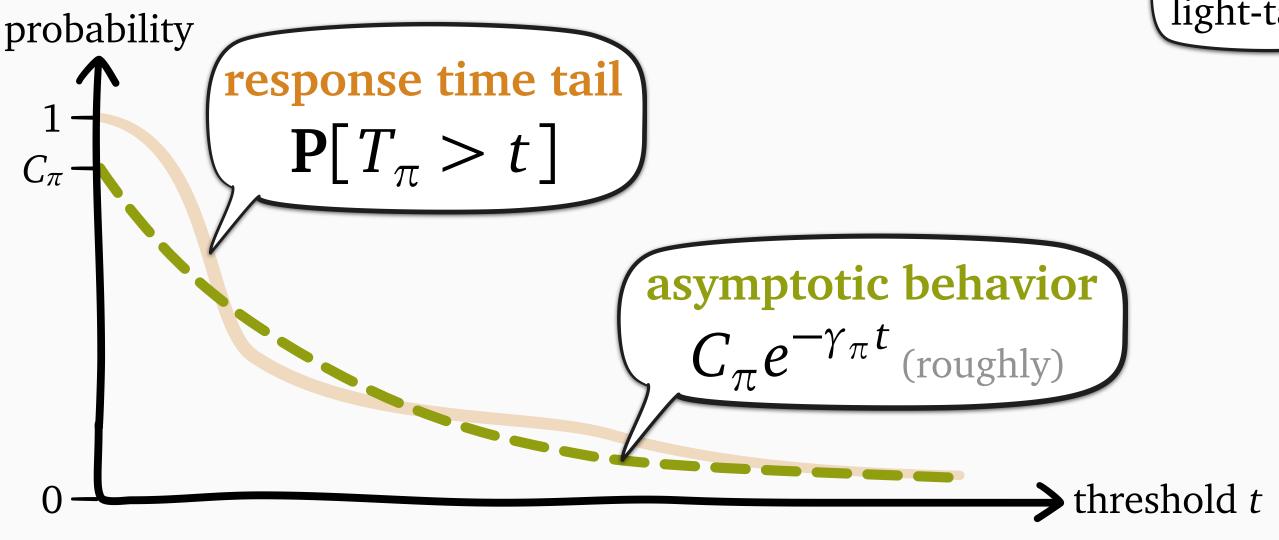
$$\gamma_{\pi} = decay \ rate \ of \ \pi$$

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



## Asymptotic response time taily

when *S* is light-tailed



Weak optimality:  $\leftarrow$  optimal  $\gamma_{\pi}$ 

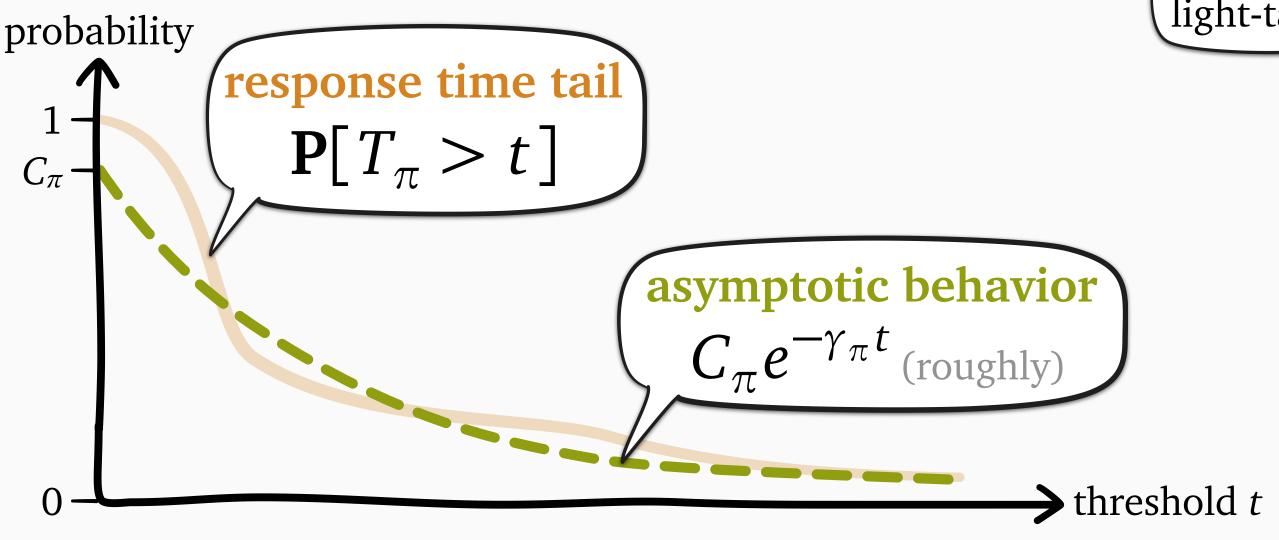
$$\gamma_{\pi} = decay \ rate \ of \ \pi$$
 $C_{\pi} = tail \ constant \ of \ \pi$ 

Strong optimality: optimal  $\gamma_{\pi}$  and  $C_{\pi}$ 

depends on policy  $\pi$ 

## Asymptotic response time taily

when S is light-tailed



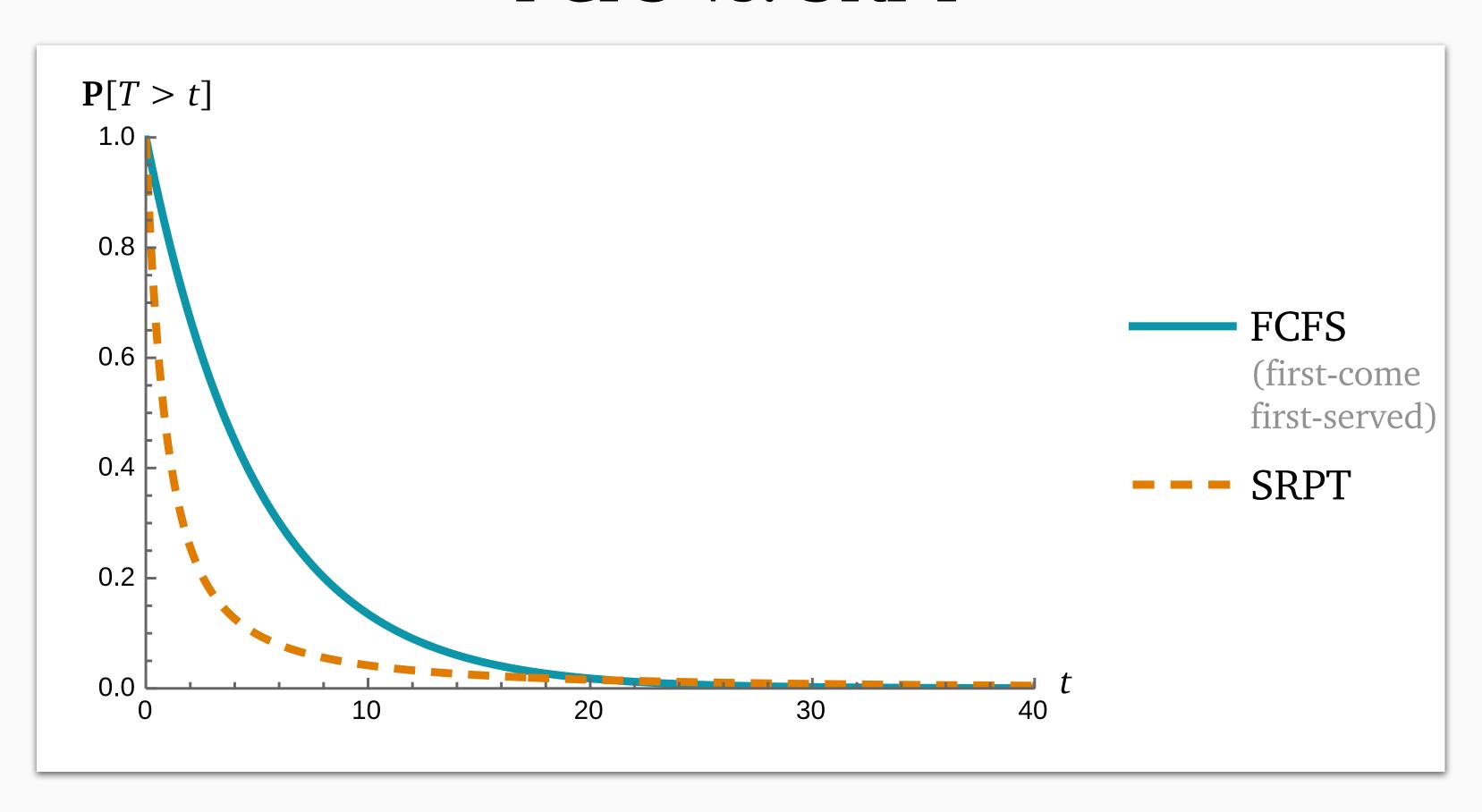
Weak optimality:  $\leftarrow$  optimal  $\gamma_{\pi}$ 

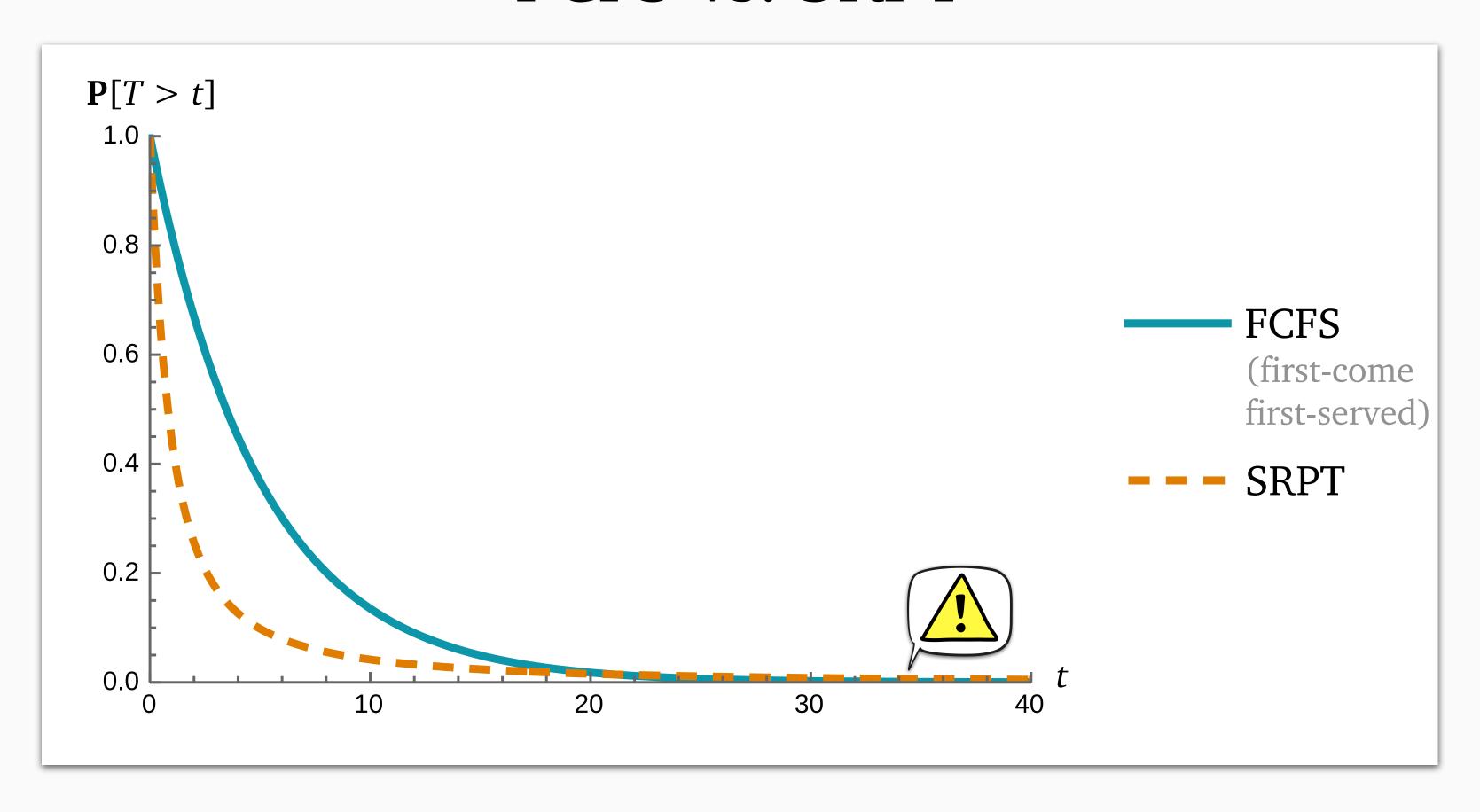
$$\gamma_{\pi} = decay \ rate \ of \ \pi$$

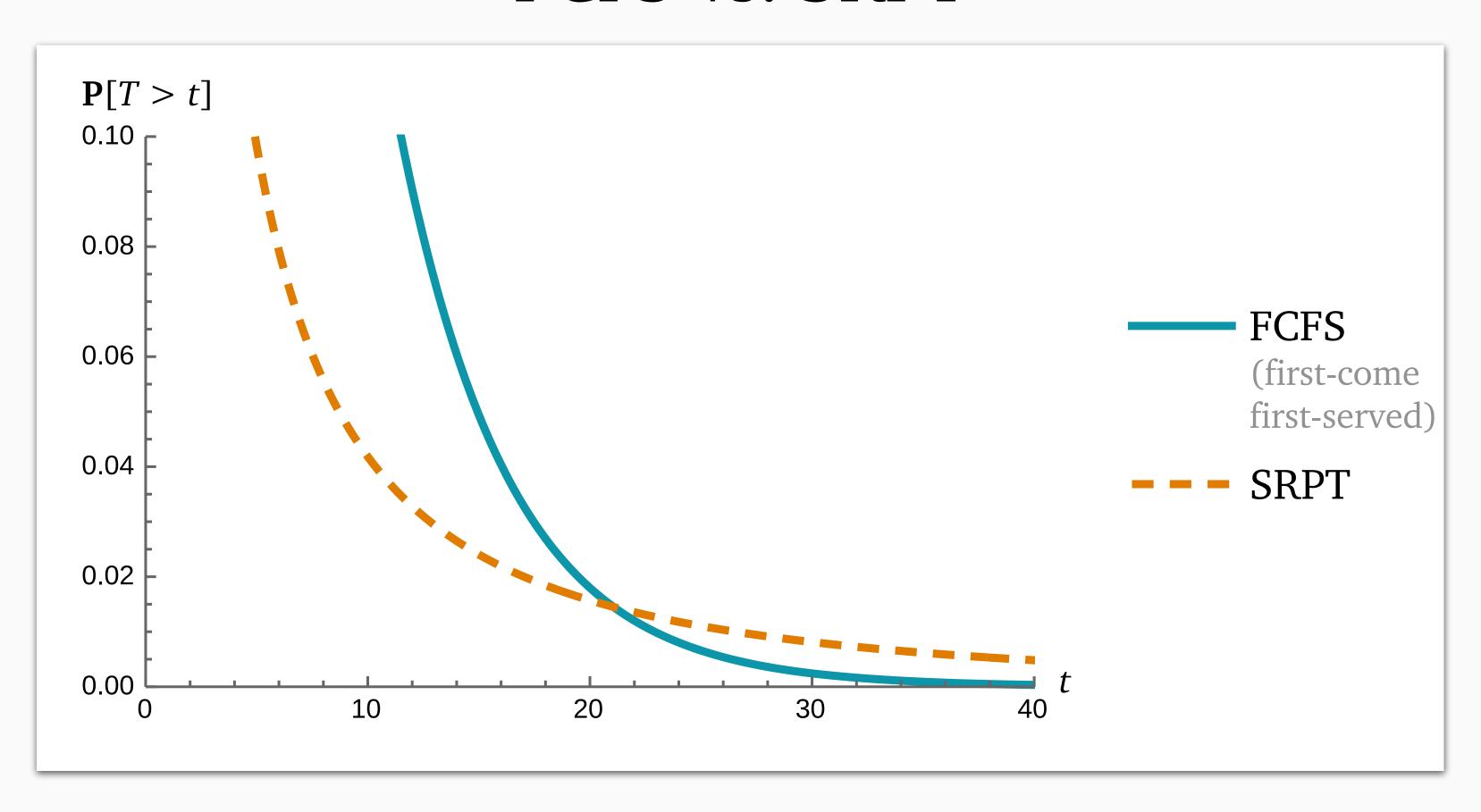
$$C_{\pi} = tail \ constant \ of \ \pi$$

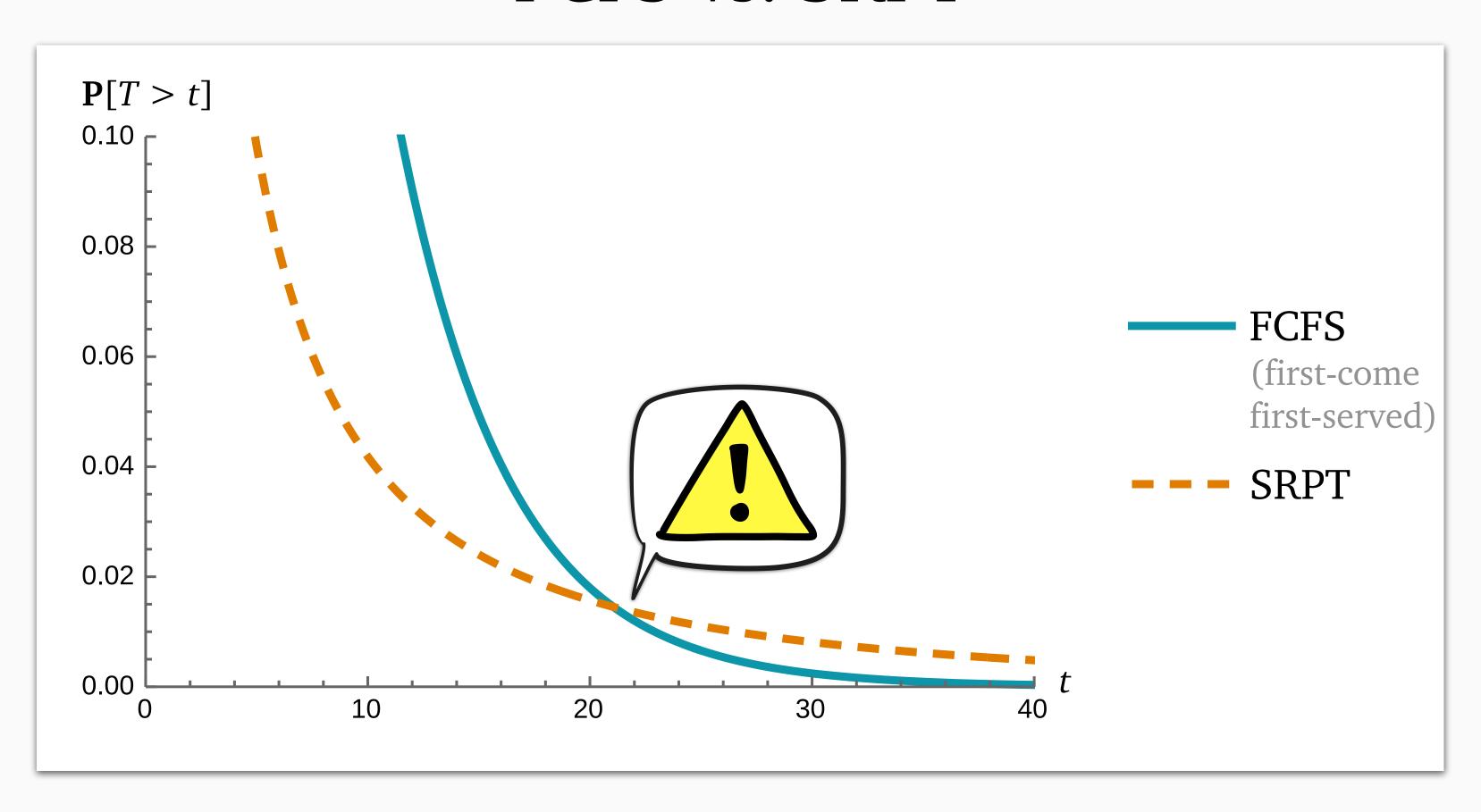
Strong optimality: optimal  $\gamma_{\pi}$  and  $C_{\pi}$ 

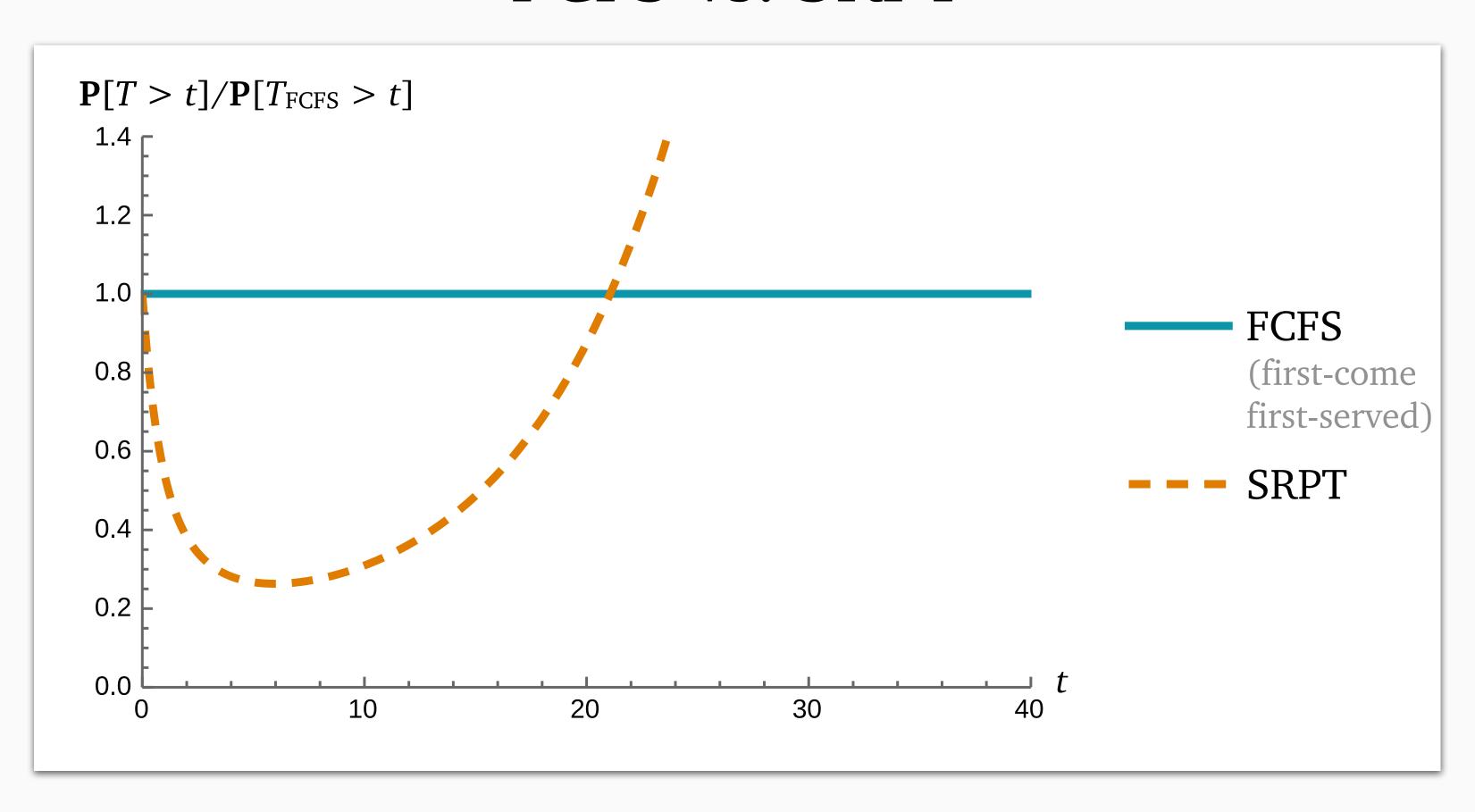
(roughly)

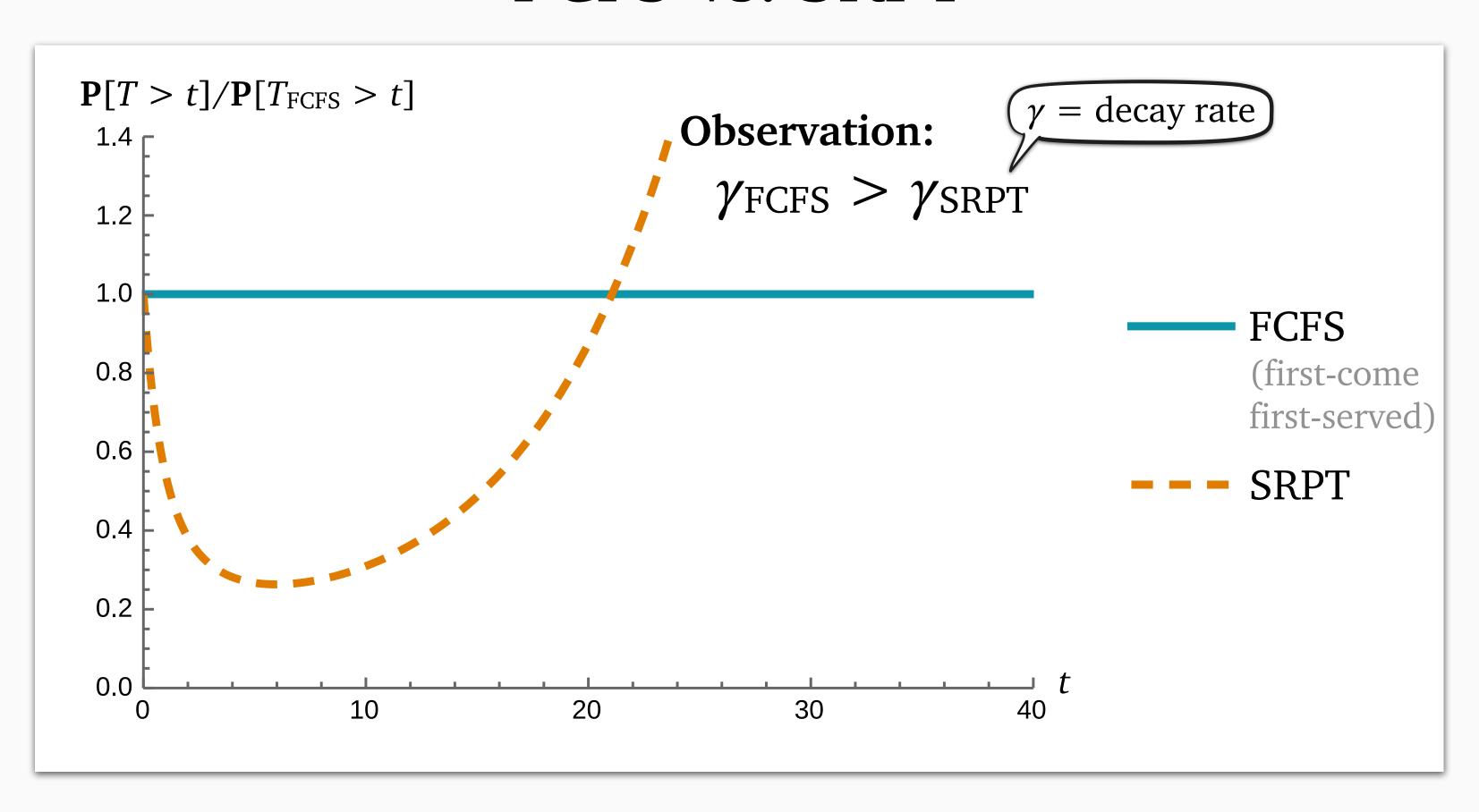


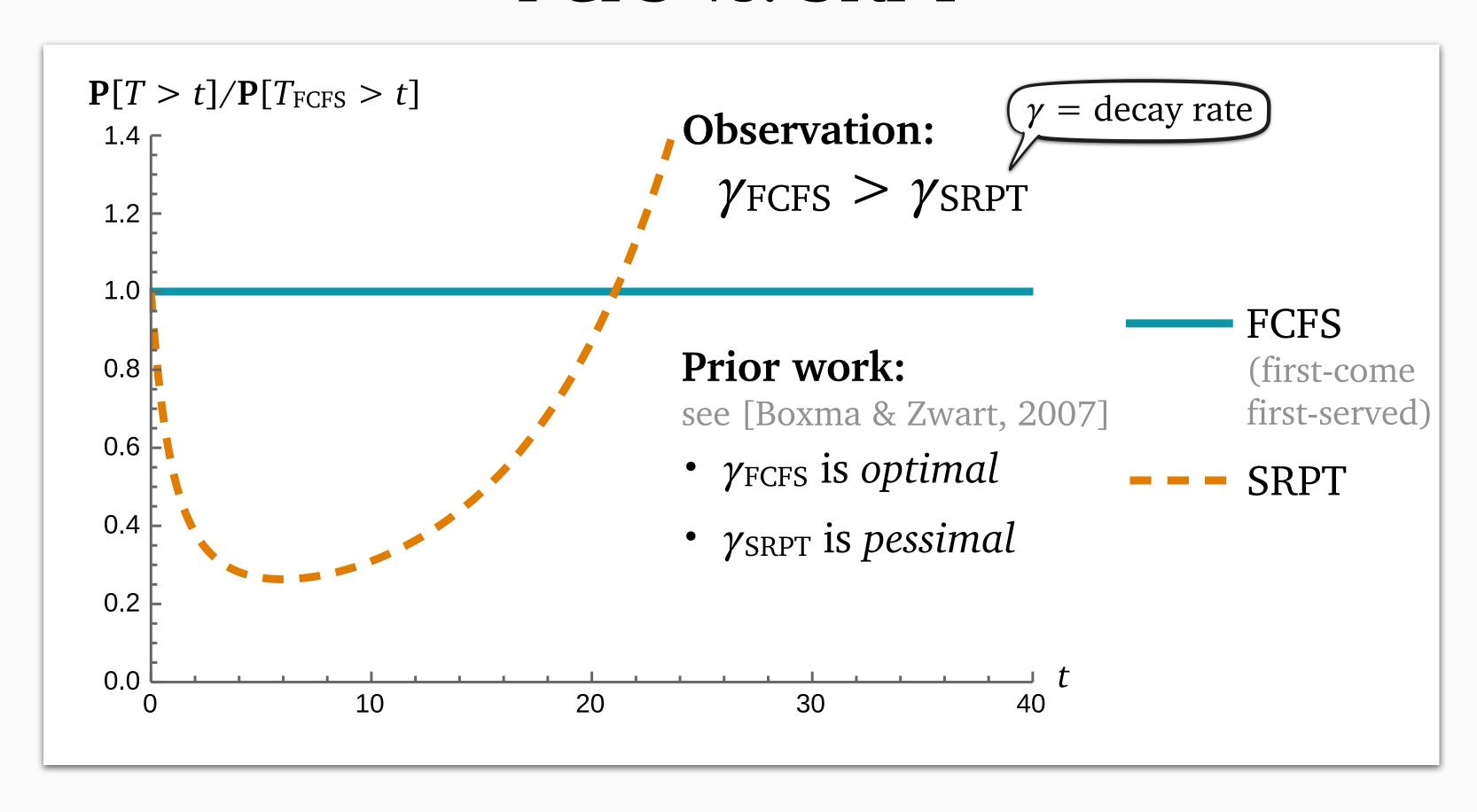


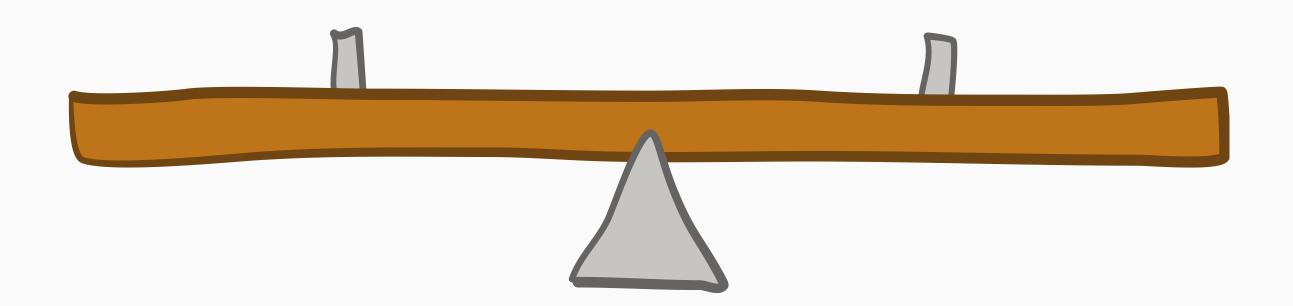


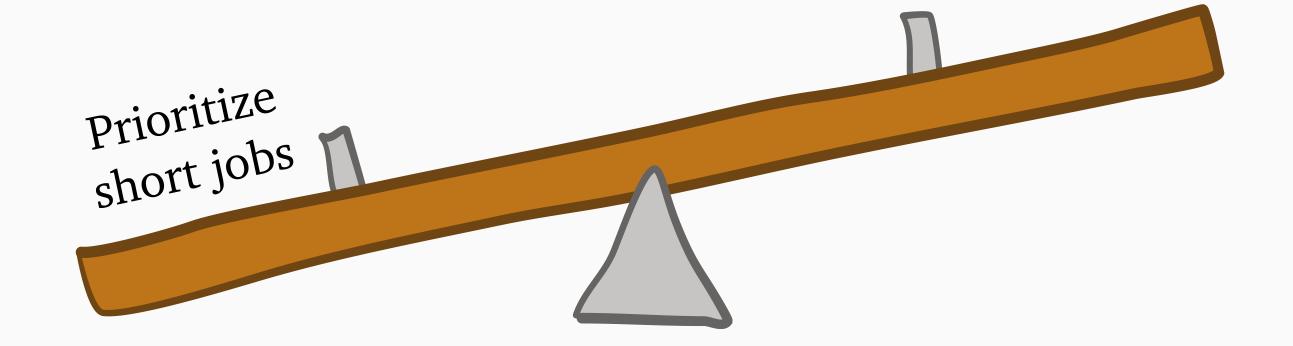


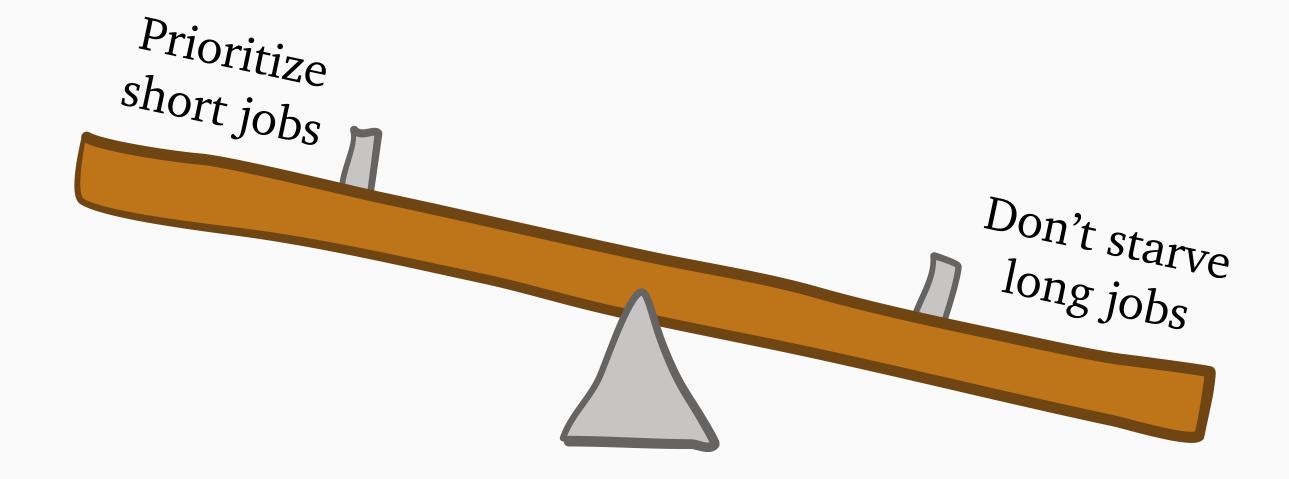


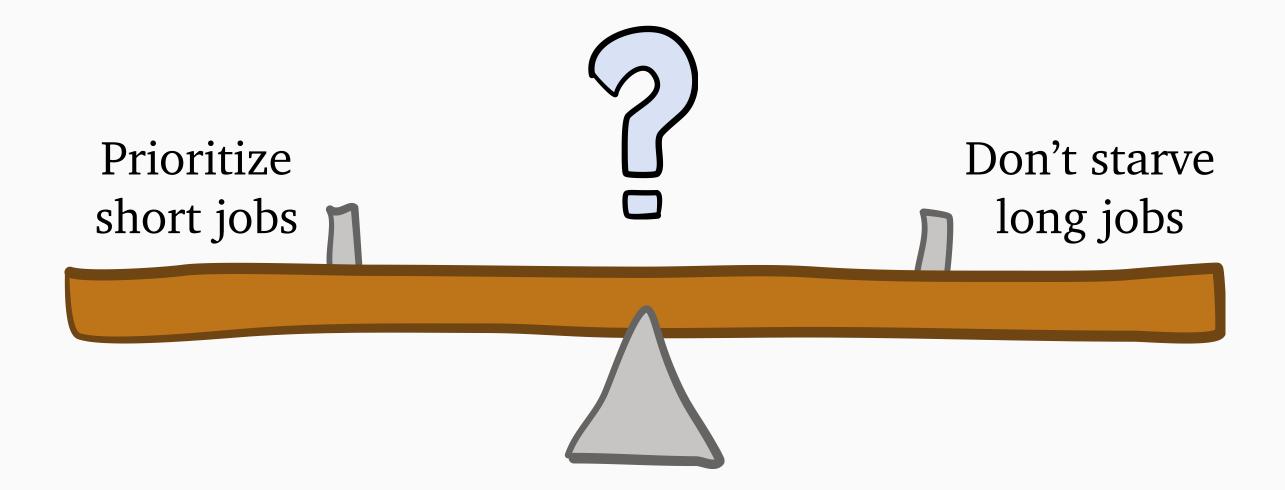


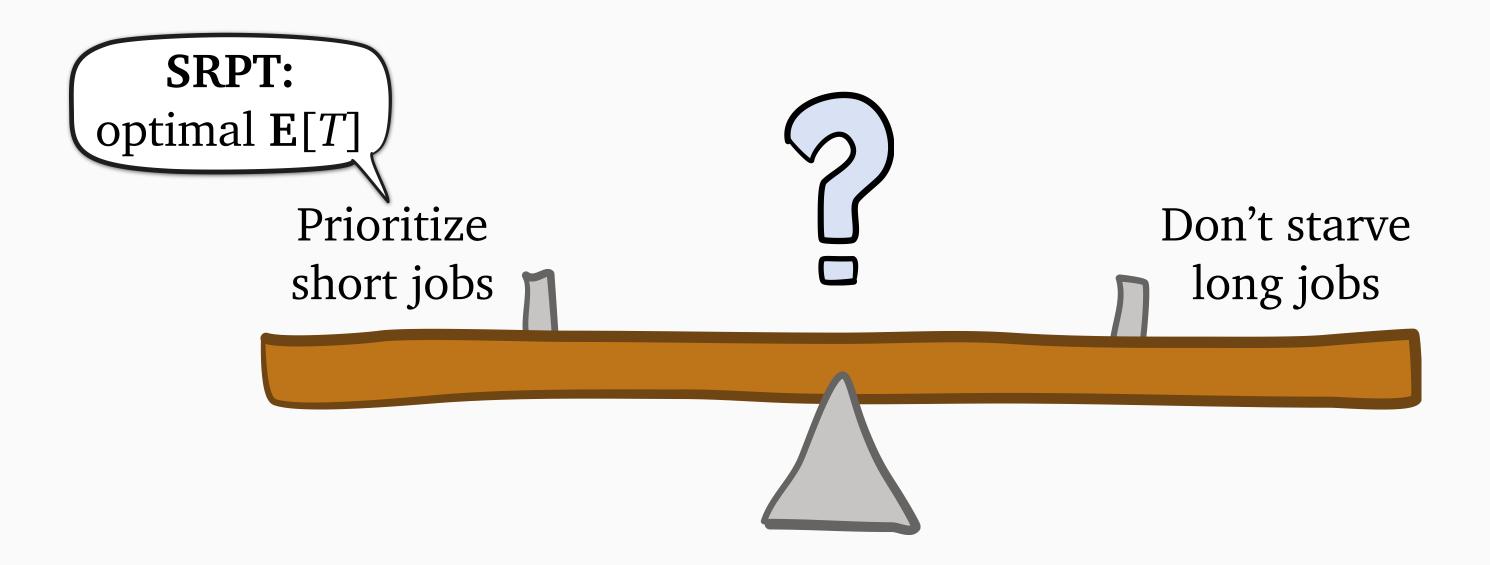


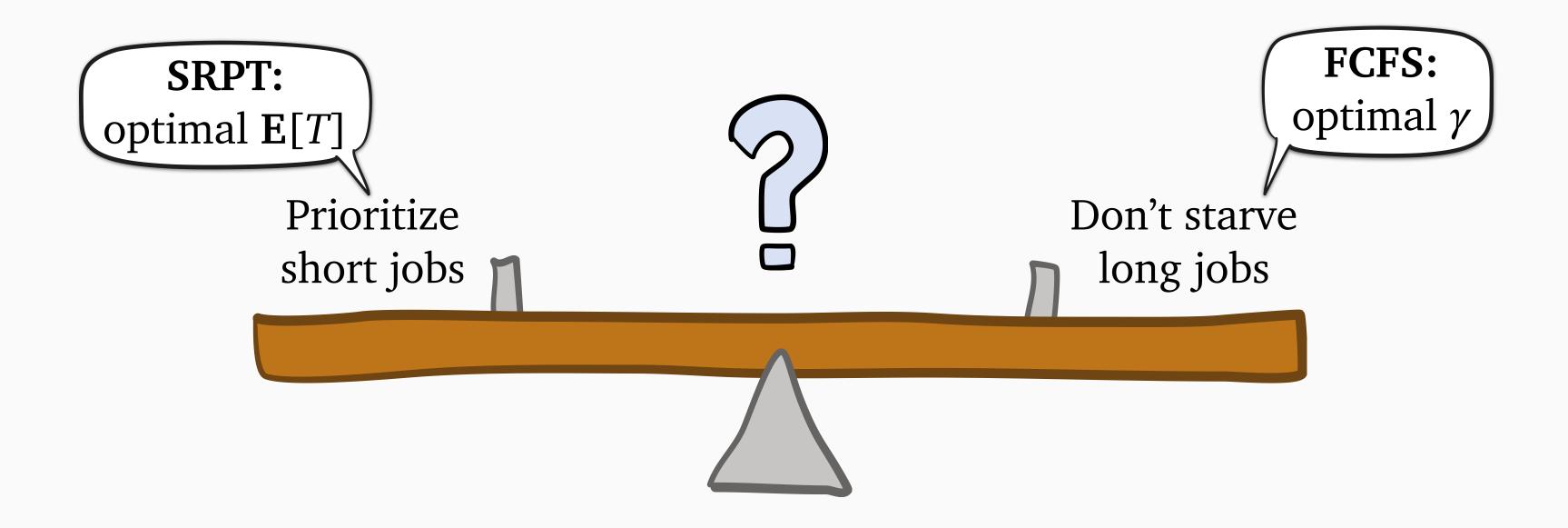


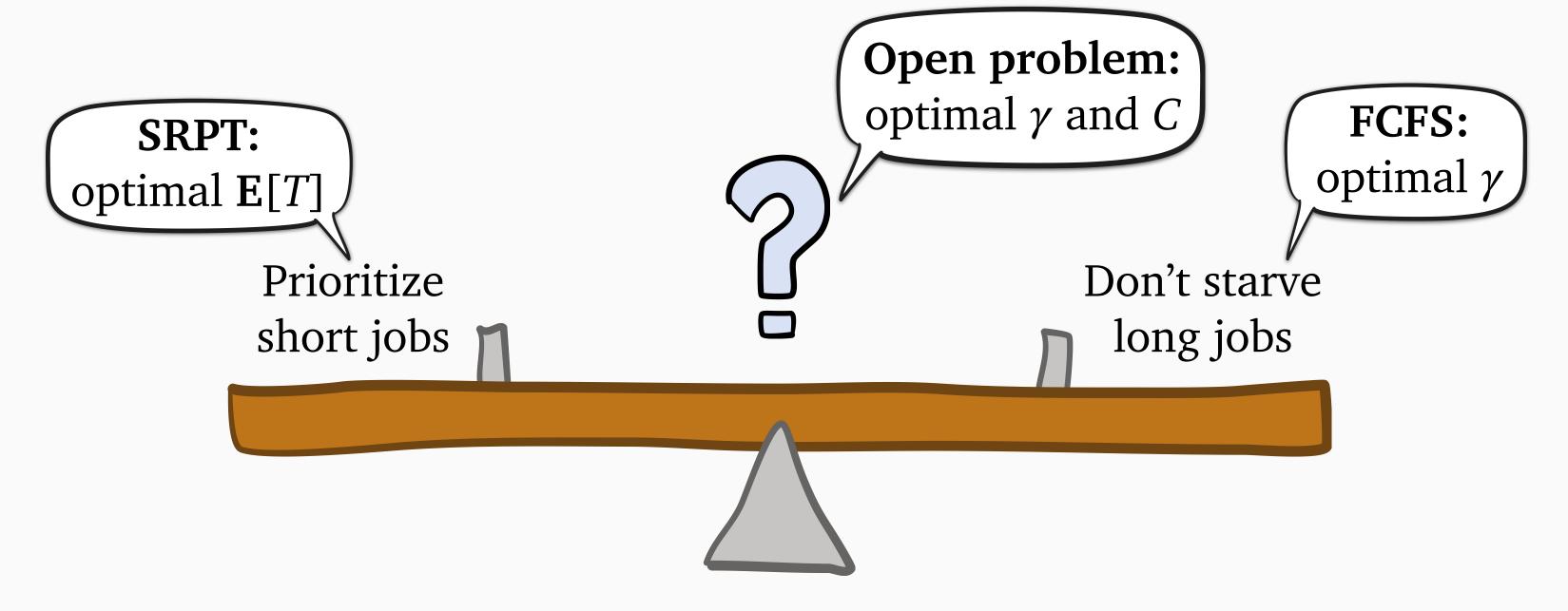


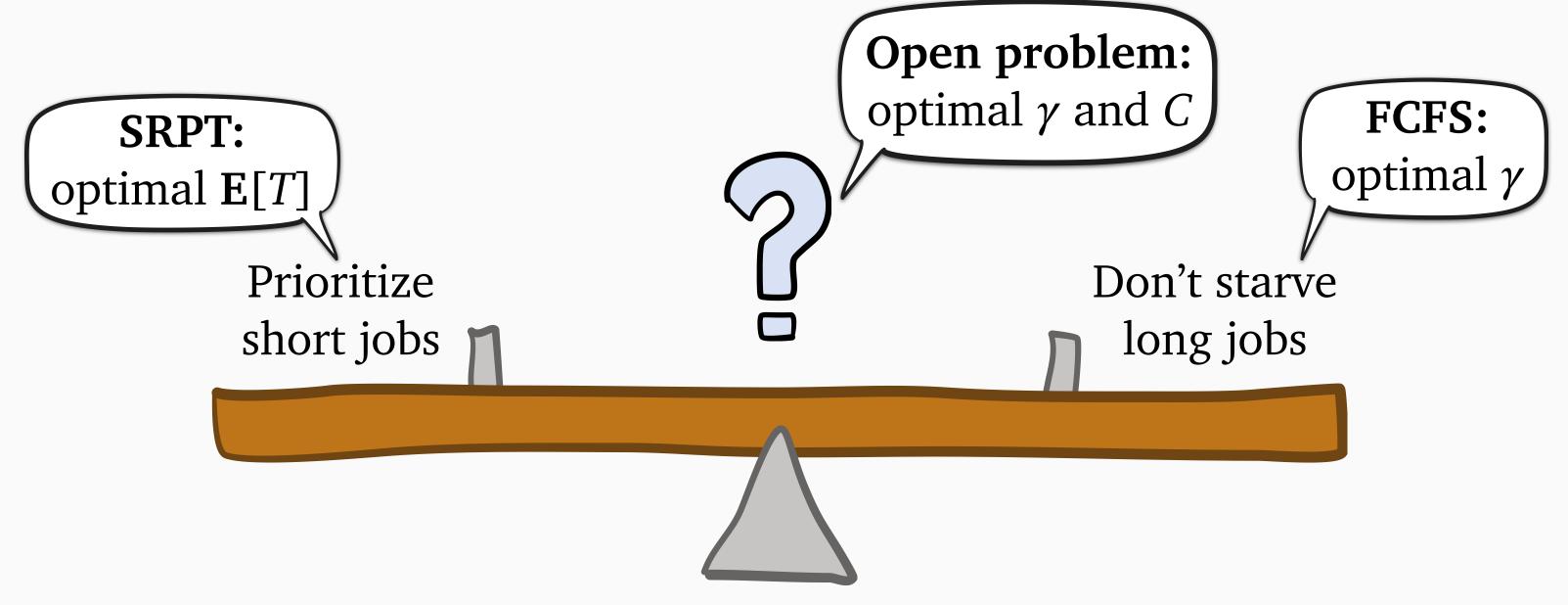






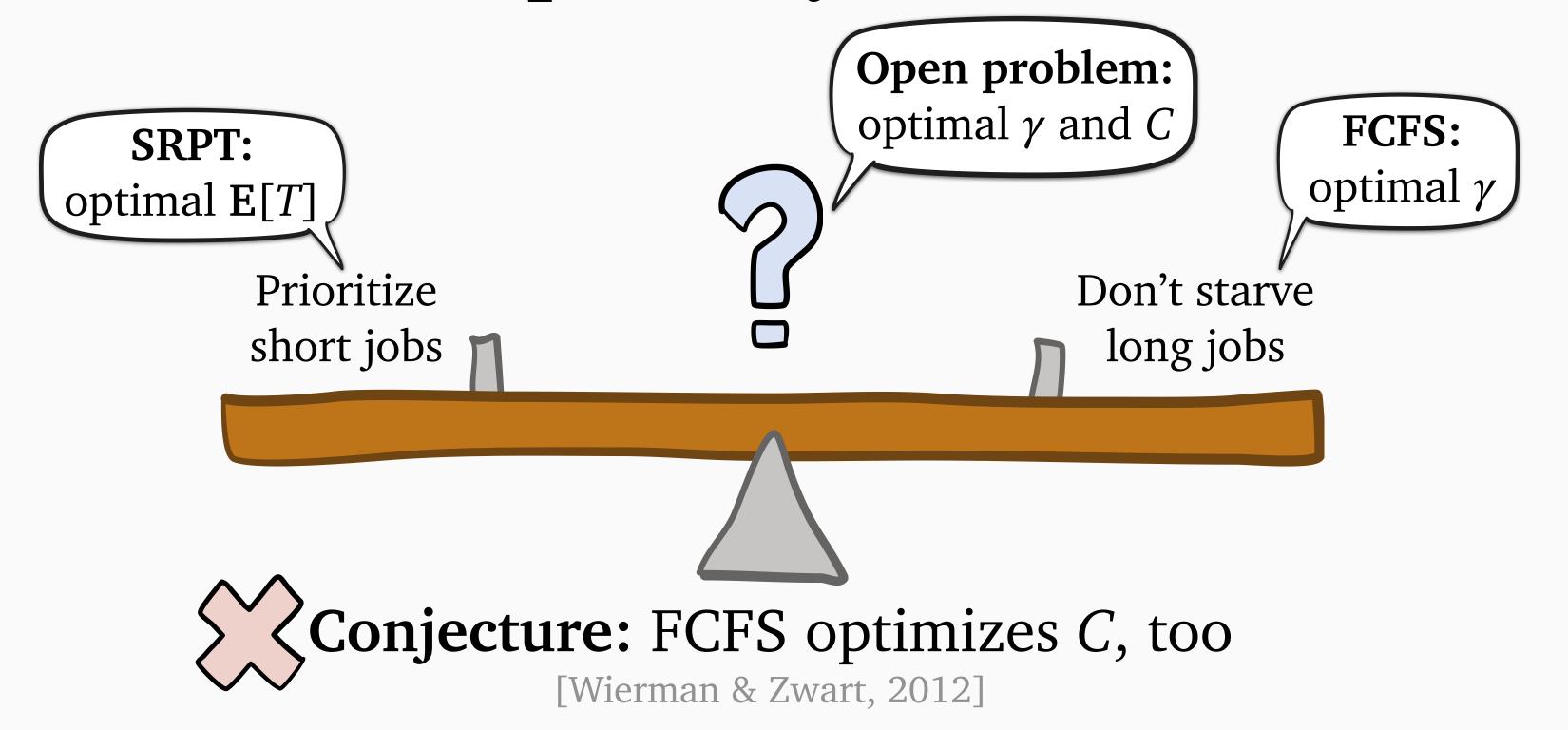


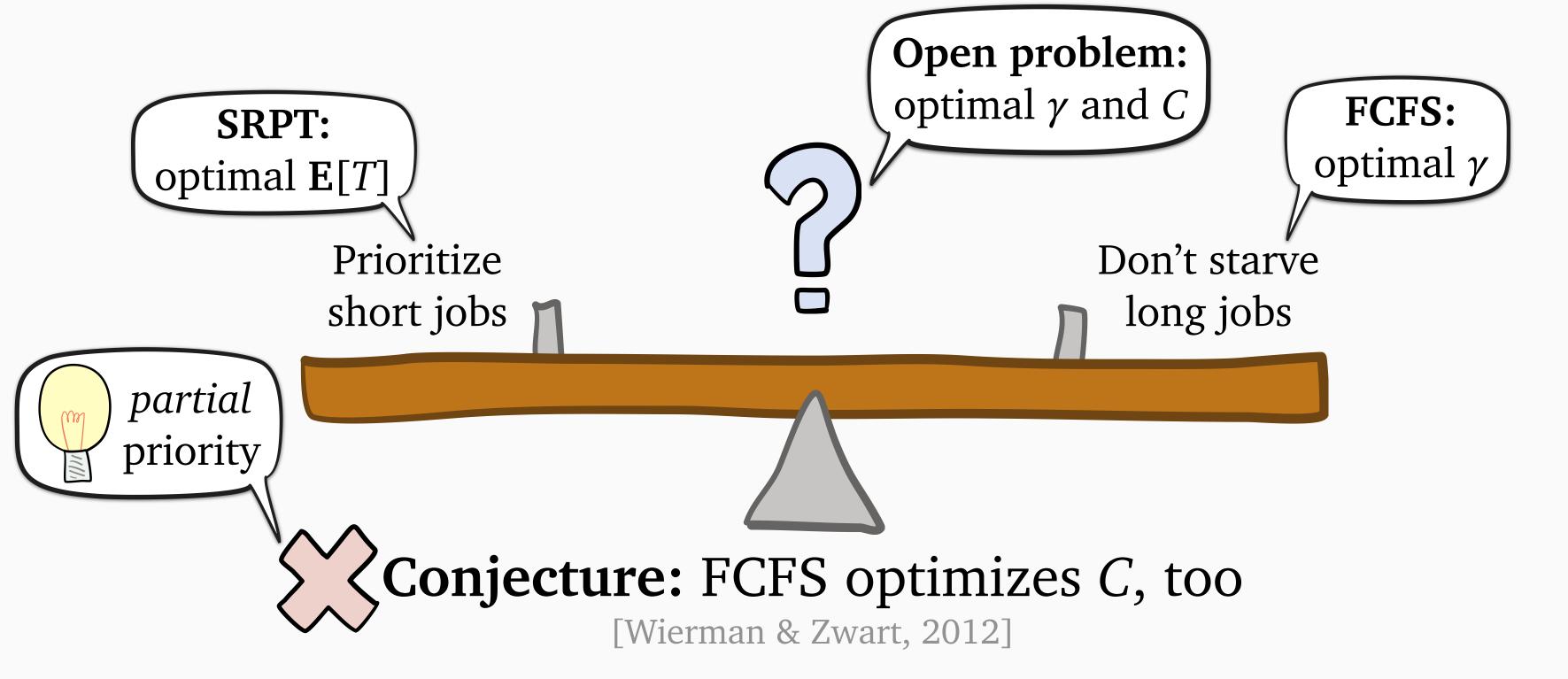


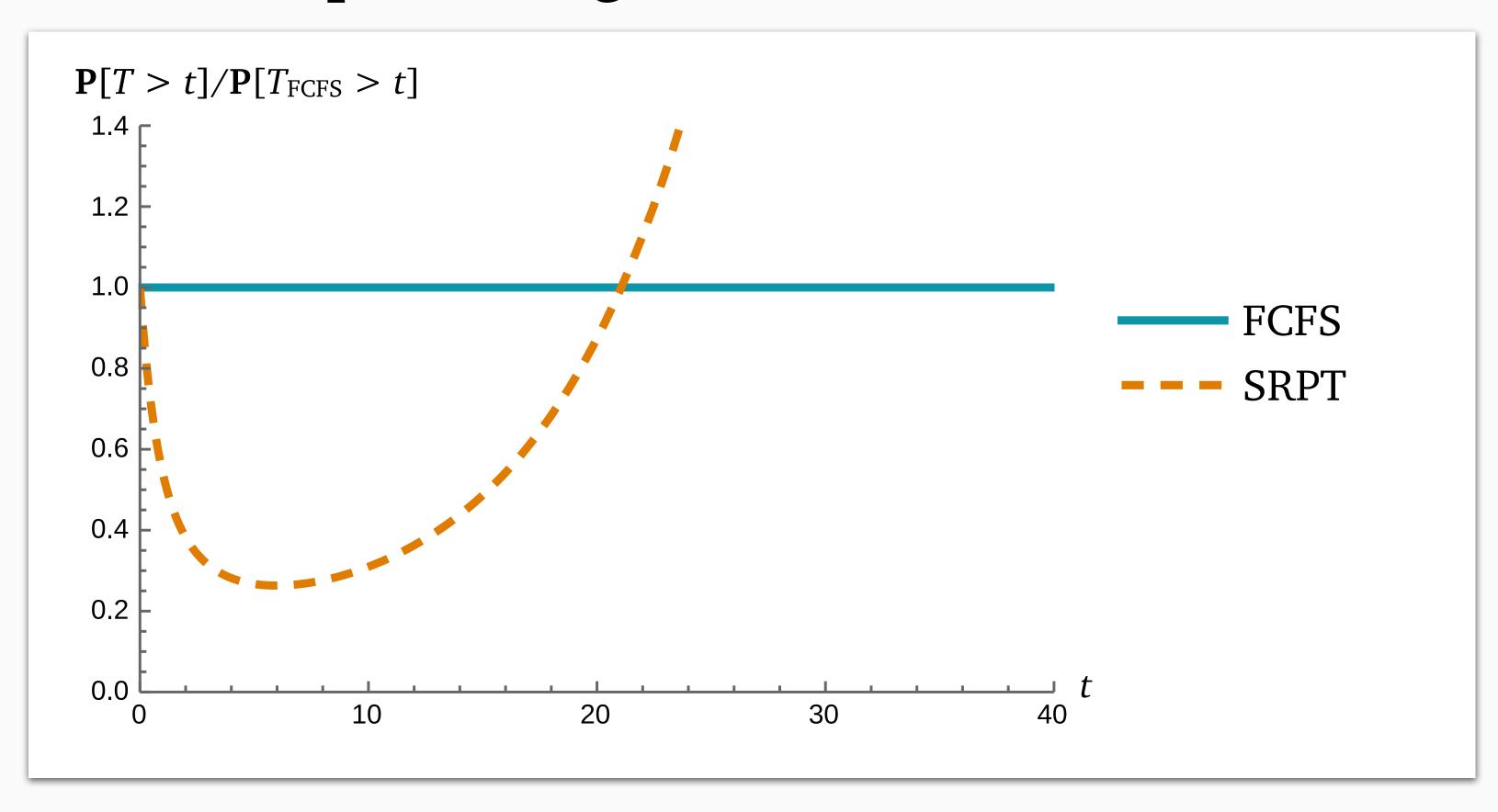


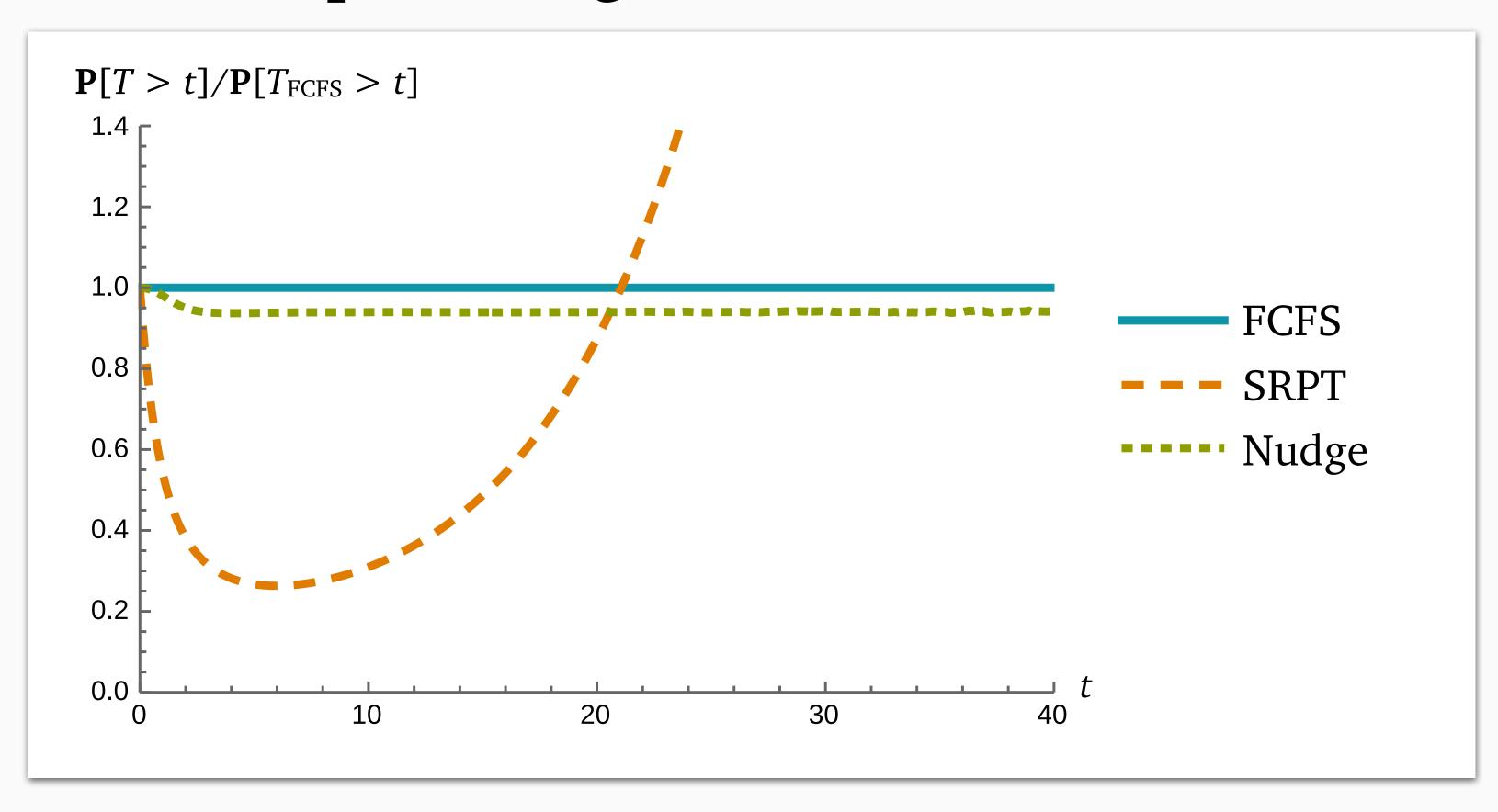
Conjecture: FCFS optimizes C, too

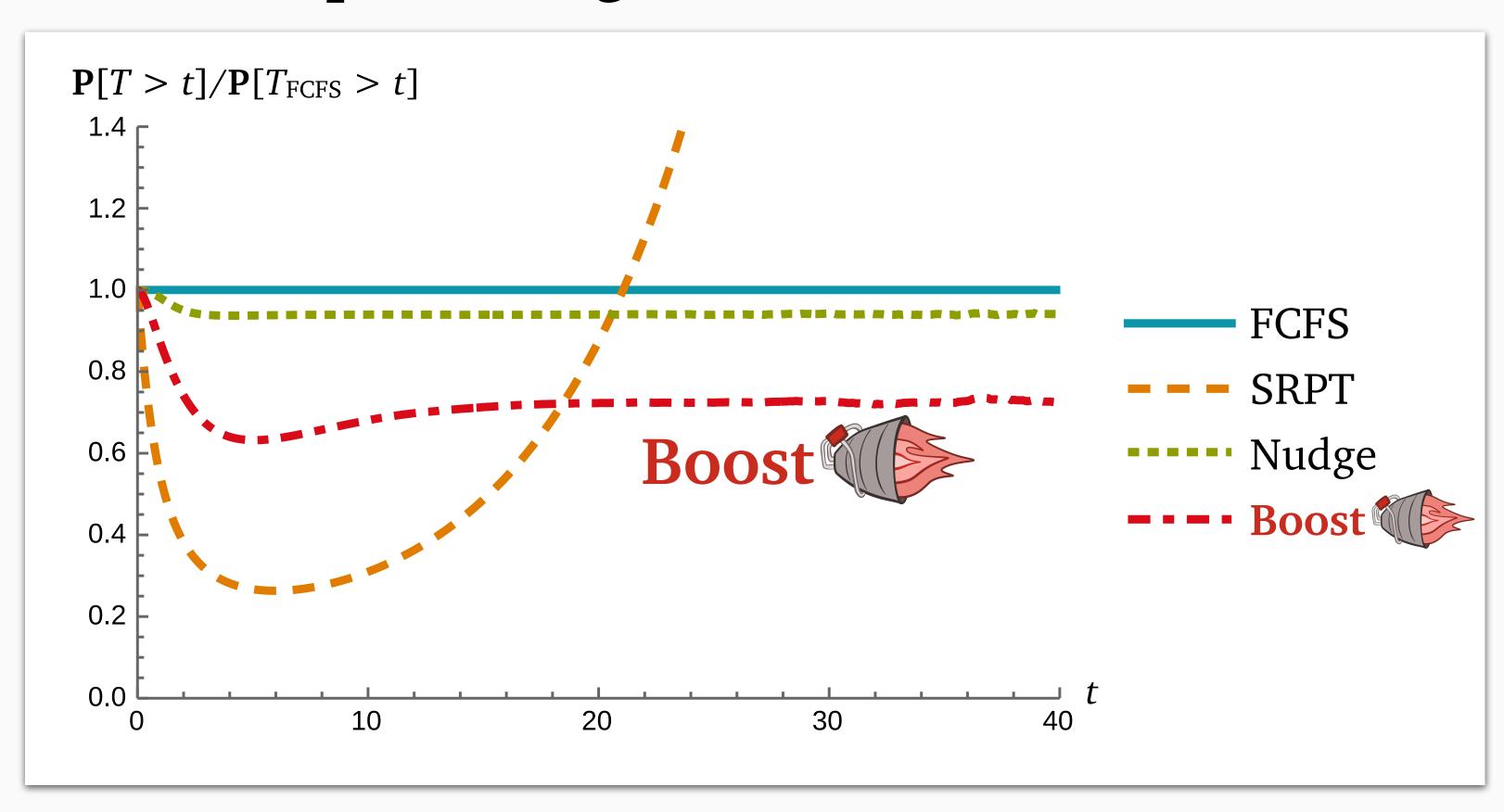
[Wierman & Zwart, 2012]

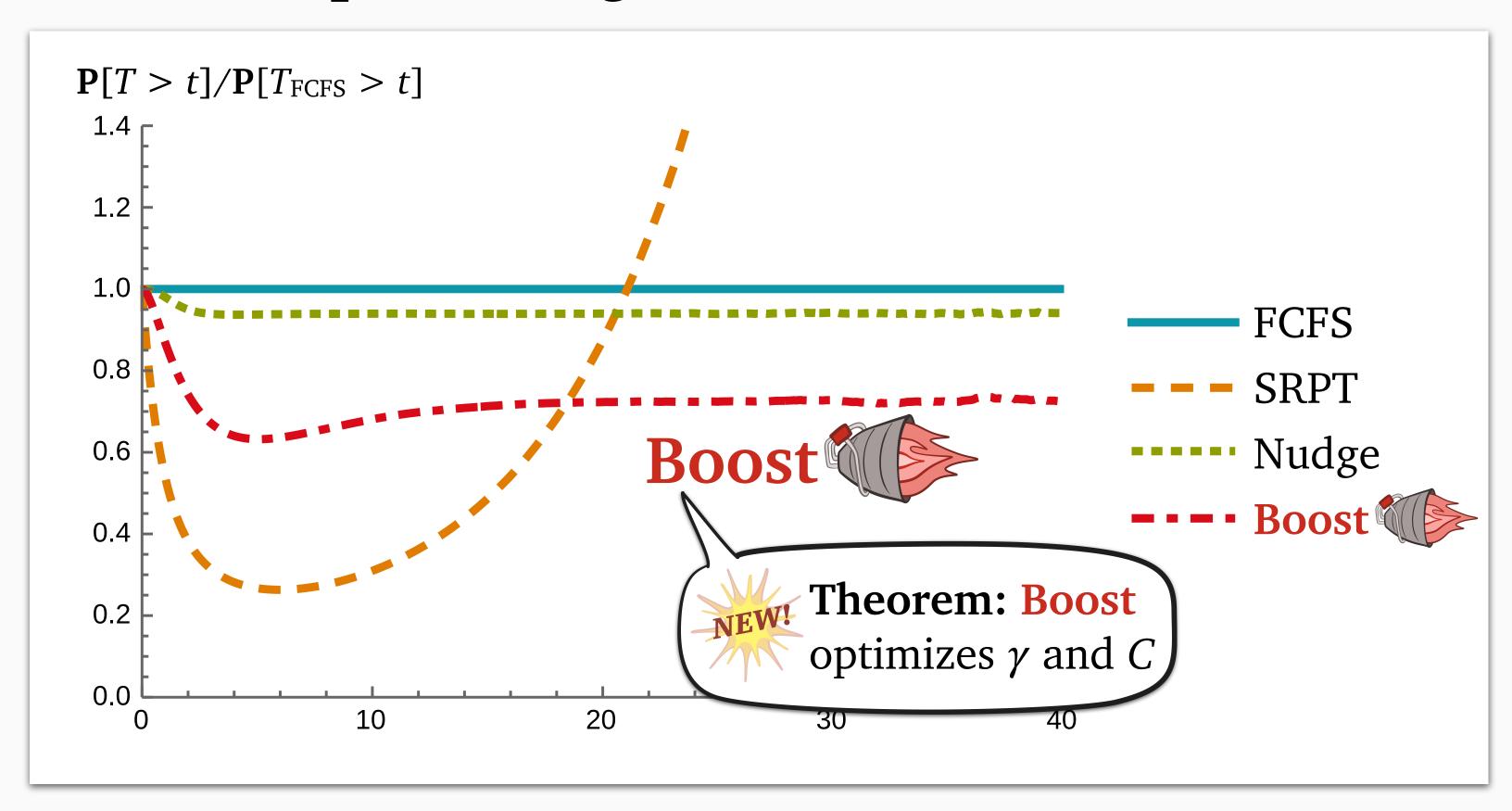














Design the **Boost** scheduling policy



Analyze **Boost**'s performance



actually a family of many policies



Design the **Boost** scheduling policy



Analyze **Boost**'s performance



actually a family of many policies

all instances



Design the **Boost** scheduling policy



Analyze **Boost**'s performance



actually a family of many policies



Design the Boost scheduling policy



Analyze **Boost**'s performance

all instances specific instance

called **γ-Boost** 



actually a family of many policies



Design the **Boost** scheduling policy



Analyze **Boost**'s performance

all instances

specific instance called γ-Boost



Prove Boost is strongly tail-optimal for light-tailed sizes

#### Known job sizes

Yu & Scully. Strongly Tail-Optimal Scheduling in the Light-Tailed *M/G/1*. SIGMETRICS 2024.

actually a family of many policies



Design the **Boost** scheduling policy



Analyze **Boost**'s performance

all instances

specific instance called γ-Boost



Prove Boost is strongly tail-optimal for light-tailed sizes



#### Known job sizes

Yu & Scully. Strongly Tail-Optimal Scheduling in the Light-Tailed M/G/1. SIGMETRICS 2024.

#### Unknown job sizes

Harley, Yu, & Scully. A Gittins Policy for Optimizing Tail Latency. SIGMETRICS 2025.

actually a family of many policies



Design the **Boost** scheduling policy



Analyze **Boost**'s performance

all instances

specific instance called γ-Boost



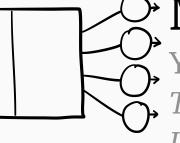
Prove Boost is strongly tail-optimal for light-tailed sizes

#### Known job sizes

Yu & Scully. Strongly Tail-Optimal Scheduling in the Light-Tailed M/G/1. SIGMETRICS 2024.

#### Unknown job sizes

Harlev, Yu, & Scully. A Gittins Policy for Optimizing Tail Latency. SIGMETRICS 2025.



#### Multiple servers

Yu, Harlev, Adakroy, & Scully. A Tale of Two Traffics: Optimizing Tail Latency in the Light-Tailed M/G/k. SIGMETRICS 2026.





Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



How do we achieve strong tail optimality?





Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



How do we achieve strong tail optimality?

# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



How do we achieve strong tail optimality?

"S Pareto-ish" (regularly varying)

$$\mathbf{P}[S > s] \sim As^{-\alpha}$$

#### Light-tailed sizes

$$P[S > s] \sim Ae^{-\alpha s}$$

"S Pareto-ish" (regularly varying)

$$P[S > s] \sim As^{-\alpha}$$

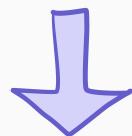
#### Light-tailed sizes

$$\mathbf{P}[S > s] \sim Ae^{-\alpha s}$$

$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} e^{-\gamma_{\pi} t}$$

"S Pareto-ish" (regularly varying)

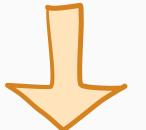
$$P[S > s] \sim As^{-\alpha}$$



$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} t^{-\gamma_{\pi}}$$

#### Light-tailed sizes

$$\mathbf{P}[S > s] \sim Ae^{-\alpha s}$$



$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} e^{-\gamma_{\pi} t}$$

"S Pareto-ish" (regularly varying)

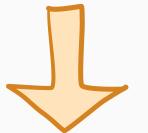
$$\mathbf{P}[S > s] \sim As^{-\alpha}$$



$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} t^{-\gamma_{\pi}}$$

## Light-tailed sizes

$$P[S > s] \sim Ae^{-\alpha s}$$



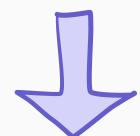
$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} e^{-\gamma_{\pi} t}$$

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

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

"S Pareto-ish" (regularly varying)

$$\mathbf{P}[S > s] \sim As^{-\alpha}$$

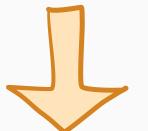


$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} t^{-\gamma_{\pi}}$$

## Light-tailed sizes

"S exponential-ish or lighter" (class I)

$$P[S > s] \sim Ae^{-\alpha s}$$



$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} e^{-\gamma_{\pi} t}$$

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

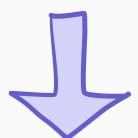
$$C_{\pi}$$
 = tail constant of  $\pi$ 

#### Weak optimality:

maximize  $\gamma_{\pi}$ 

"S Pareto-ish" (regularly varying)

$$\mathbf{P}[S > s] \sim As^{-\alpha}$$

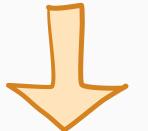


$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} t^{-\gamma_{\pi}}$$

## Light-tailed sizes

"S exponential-ish or lighter" (class I)

$$P[S > s] \sim Ae^{-\alpha s}$$



$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} e^{-\gamma_{\pi} t}$$

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

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

Weak optimality:

maximize  $\gamma_{\pi}$ 

Strong optimality:

maximize  $\gamma_{\pi}$ , minimize  $C_{\pi}$ 

"S Pareto-ish" (regularly varying)

$$\mathbf{P}[S > s] \sim As^{-\alpha}$$

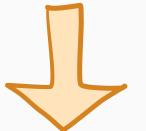


$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} t^{-\gamma_{\pi}}$$

(roughly)

"S exponential-ish or lighter" (class I)

$$P[S > s] \sim Ae^{-\alpha s}$$



$$\mathbf{P}[T_{\pi} > t] \sim C_{\pi} e^{-\gamma_{\pi} t}$$

(roughly)

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

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

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

#### Weak optimality:

maximize  $\gamma_{\pi}$ 

#### Strong optimality:

maximize  $\gamma_{\pi}$ , minimize  $C_{\pi}$ 

(roughly)

Heavy-tailed sizes	Light-tailed sizes

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)		
FCFS		

<b>T T</b>	. • 1	1	•
$H_{\Omega} \chi_{\Lambda} \chi_{\Lambda}$	1911		C17AC
Heavy-1	Lall	CU	<b>217C2</b>

Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

optimal  $\gamma = \alpha$ 

**FCFS** 

pessimal  $\gamma = \alpha - 1$ 

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ
Main cause of large T?		

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ
Main cause of large T?	"Catastrophe" one giant job	

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$ I'm stuck behind the giant job	optimal γ
Main cause of large T?	"Catastrophe" one giant job	

Heavy-	tailed	sizes
ricavy	LUIICU	DIZCO

Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

optimal  $\gamma = \alpha$ I'm the giant job

pessimal  $\gamma$ 

**FCFS** 

Main cause of large T?

pessimal  $\gamma = \alpha - 1$ I'm stuck behind the giant job

optimal  $\gamma$ 

"Catastrophe" one giant job

Heavy-tailed sizes

Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

optimal  $\gamma = \alpha$ I'm the giant job

pessimal  $\gamma$ 

**FCFS** 

Main cause of large T?

pessimal γ = αI'm stuck behind the giant job"Catastrophe"

Prioritize
short jobs
[Scully & van Kreveld, 2025]

**Heavy-tailed sizes** 

Light-tailed sizes

SRPT, LAS, etc.

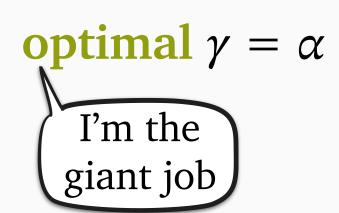
(least attain d service)

also optimal C for

heavy-tailed sizes

FCFS

Main cause of large T?



one giant job

pessimal  $\gamma$ 

pessimal  $\gamma = \alpha$ I'm stuck behind the giant job

"Catastrophe"

Prioritize short jobs short jobs [Scully & van Kreveld, 2025]

Heavy-	tailed	sizes
ricavy	LUIICU	DIZCO

Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

optimal  $\gamma = \alpha$ I'm the giant job

pessimal  $\gamma$ 

**FCFS** 

Main cause of large T?

pessimal  $\gamma = \alpha - 1$ I'm stuck behind the giant job

optimal  $\gamma$ 

"Catastrophe" one giant job

Heavy-tailed sizes
--------------------

Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

optimal  $\gamma = \alpha$ I'm the giant job

pessimal  $\gamma$ 

**FCFS** 

Main cause of large T?

pessimal  $\gamma = \alpha - 1$ I'm stuck behind the giant job

optimal  $\gamma$ 

"Catastrophe" one giant job

"Conspiracy" lots of biggish jobs

Heavy-	tail	led	sizes

Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

optimal  $\gamma = \alpha$ I'm the giant job

pessimal  $\gamma$ 

**FCFS** 

Main cause of large T?

pessimal  $\gamma = \alpha - 1$ I'm stuck behind the giant job

"Catastrophe" one giant job

Optimal  $\gamma$ I see lots of work when I arrive

"Conspiracy" lots of biggish jobs

#### Heavy-tailed sizes

#### Light-tailed sizes

SRPT, LAS, etc.

(least attained service)

**FCFS** 

Main cause of large T?

optimal  $\gamma = \alpha$ I'm the giant job

pessimal  $\gamma = \alpha - 1$ I'm stuck behind the giant job

"Catastrophe" one giant job

pessimal  $\gamma$ I'm a very big job,

lots of smaller jobs are passing me

optimal  $\gamma$ 

I see lots of *work* when I arrive

"Conspiracy" lots of biggish jobs

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ
Main cause of large T?	"Catastrophe" one giant job	"Conspiracy" lots of biggish jobs

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ
SRPT or LAS with just two buckets		
Main cause of large T?	"Catastrophe" one giant job	"Conspiracy" lots of biggish jobs

	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ
SRPT or LAS with just two buckets	pessimal $\gamma = \alpha - 1$	
Main cause of large T?	"Catastrophe" one giant job	"Conspiracy" lots of biggish jobs

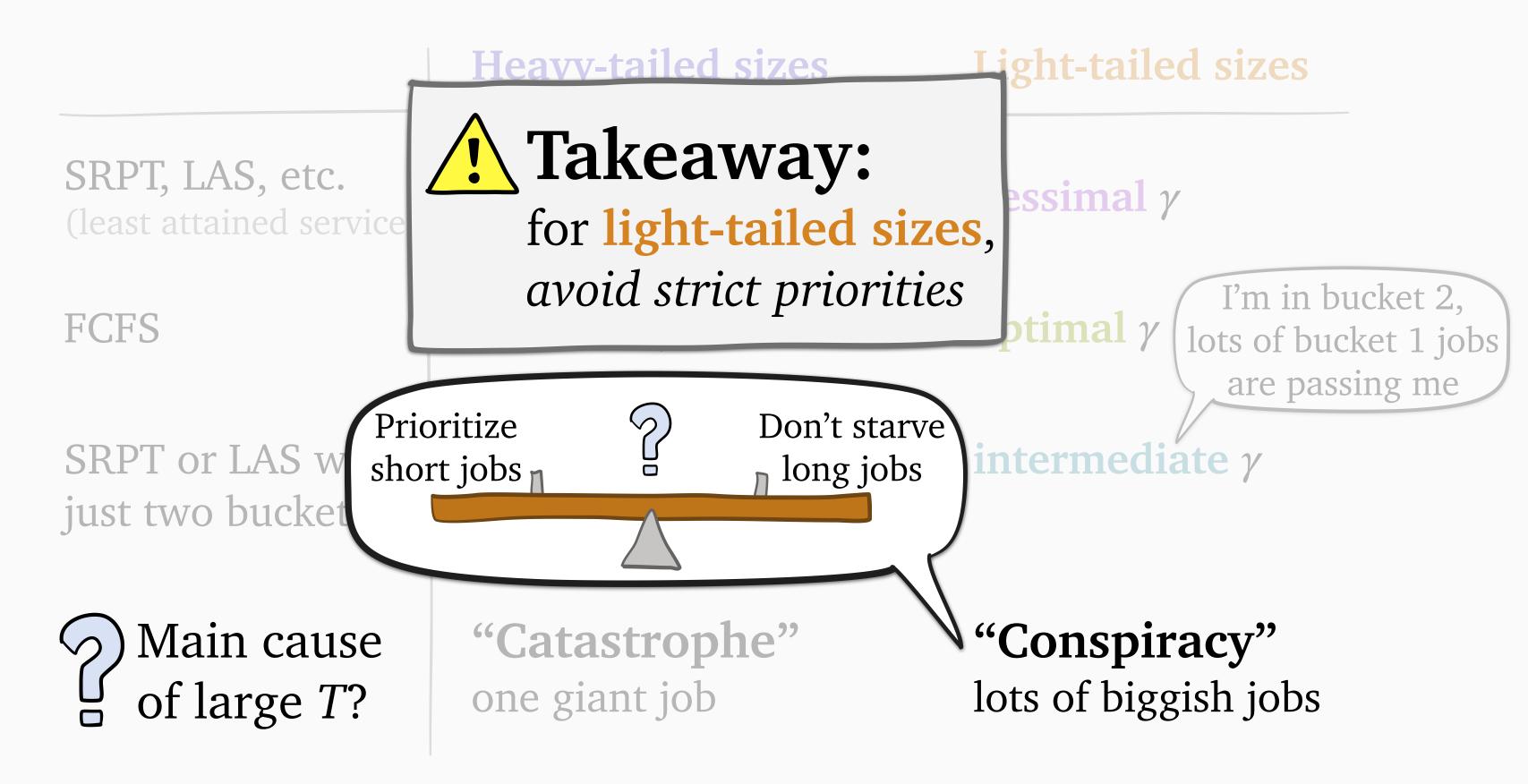
	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal $\gamma$
SRPT or LAS with just two buckets	pessimal $\gamma = \alpha - 1$	intermediate γ
Main cause of large T?	"Catastrophe" one giant job	"Conspiracy" lots of biggish jobs

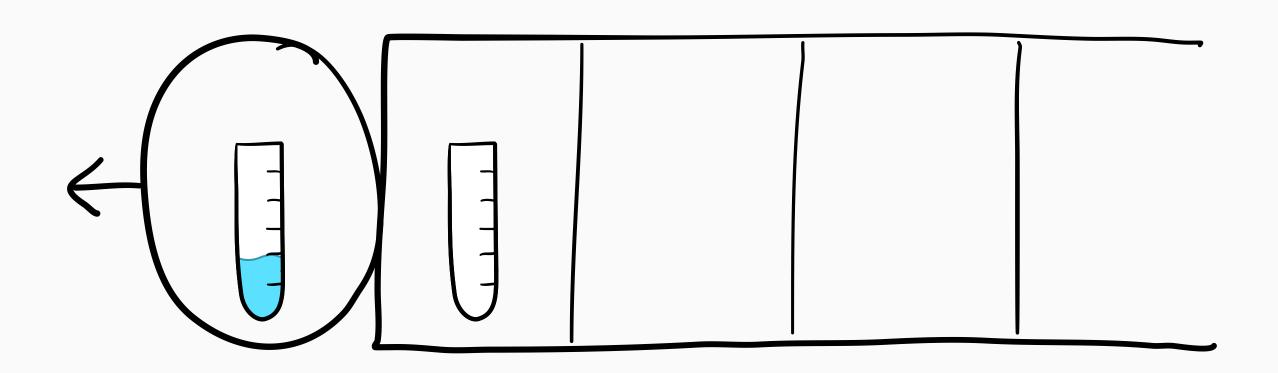
	Heavy-tailed sizes	Light-tailed sizes
SRPT, LAS, etc. (least attained service)	optimal $\gamma = \alpha$	pessimal γ
FCFS	pessimal $\gamma = \alpha - 1$	optimal γ lots of bucket 1 jobs are passing me
SRPT or LAS with just two buckets	pessimal $\gamma = \alpha - 1$	intermediate $\gamma$
Main cause of large T?	"Catastrophe" one giant job	"Conspiracy" lots of biggish jobs

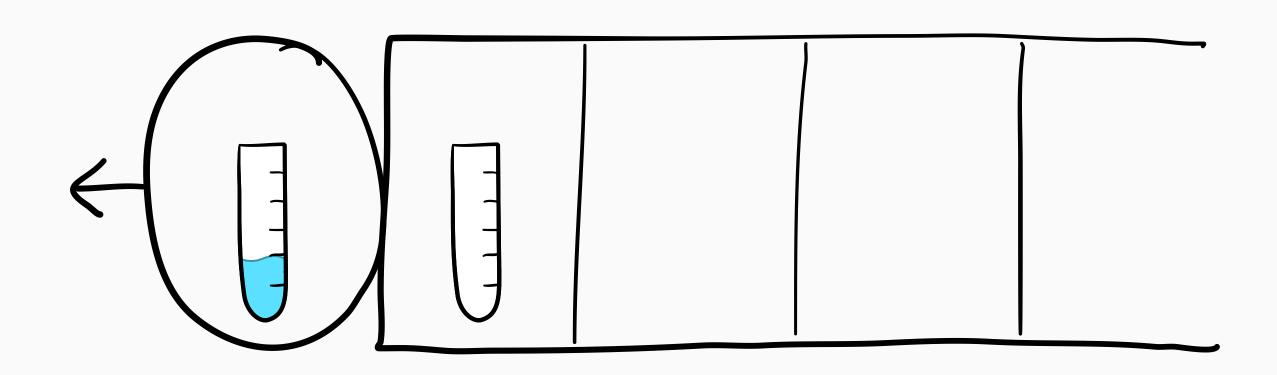
# What causes large response times?

Heavy-tailed sizes Light-tailed sizes SRPT, LAS, etc. optimal  $\gamma = \alpha$ pessimal  $\gamma$ (least attained service) I'm in bucket 2, lots of bucket 1 jobs pessimal  $\gamma = \alpha - 1$ **FCFS** are passing me Prioritize Don't starve SRPT or LAS w intermediate y short jobs n long jobs just two bucket Main cause of large *T*? "Conspiracy" "Catastrophe" lots of biggish jobs one giant job

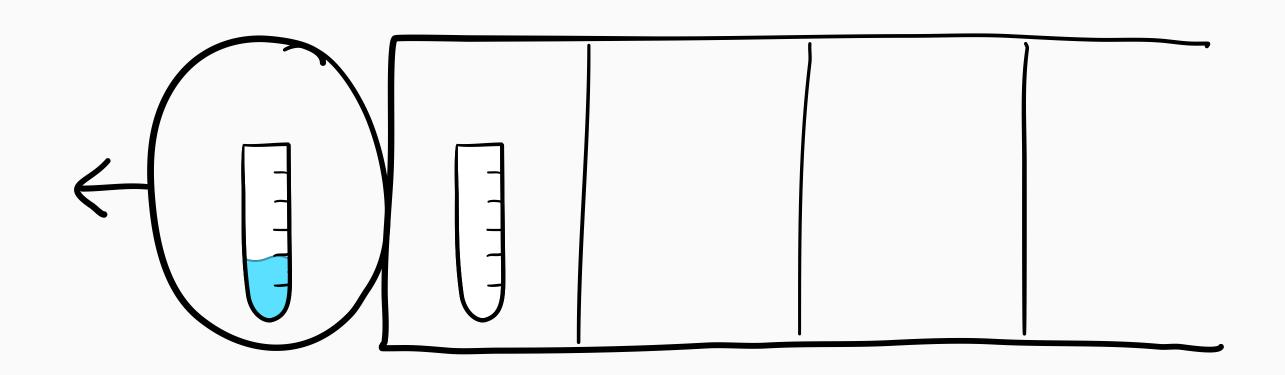
# What causes large response times?

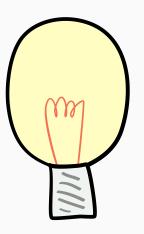




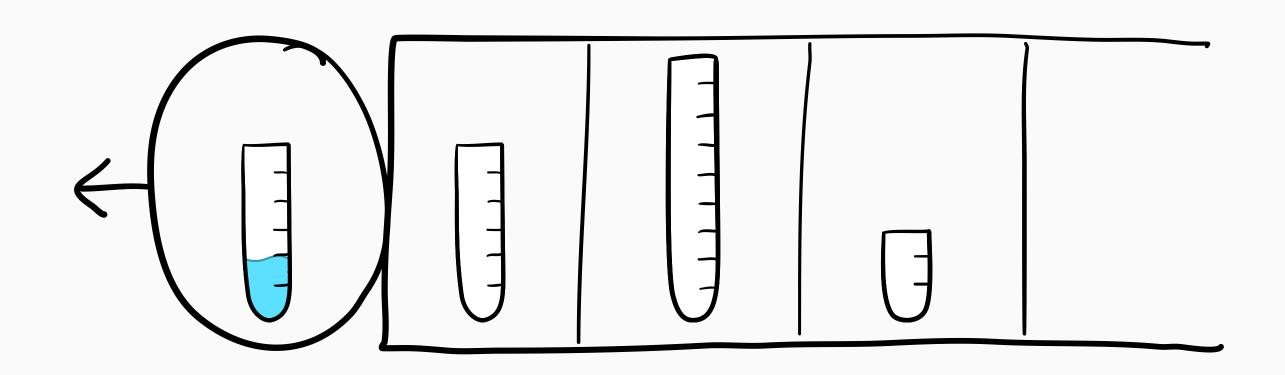


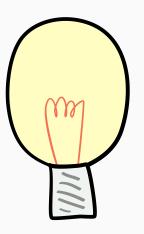




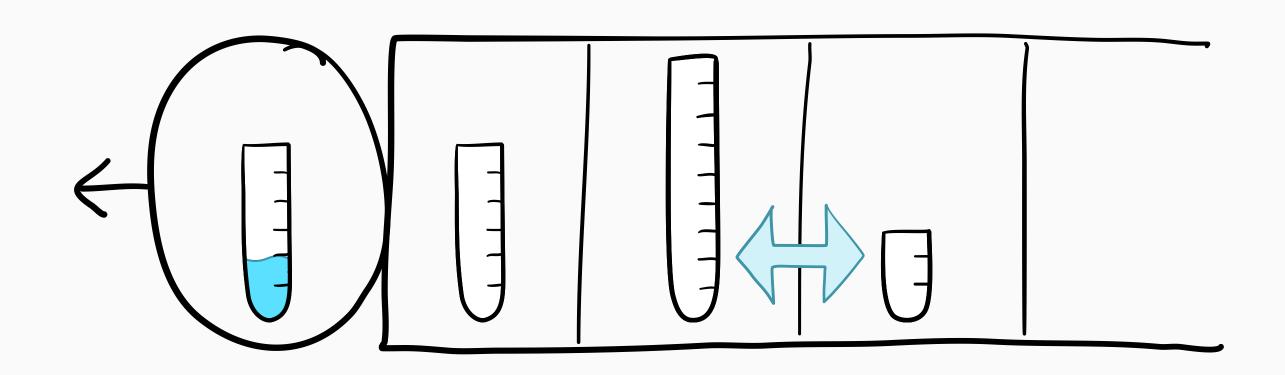


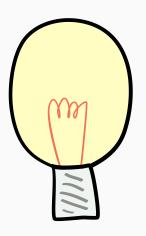
Nudge [Grosof et al., 2021]



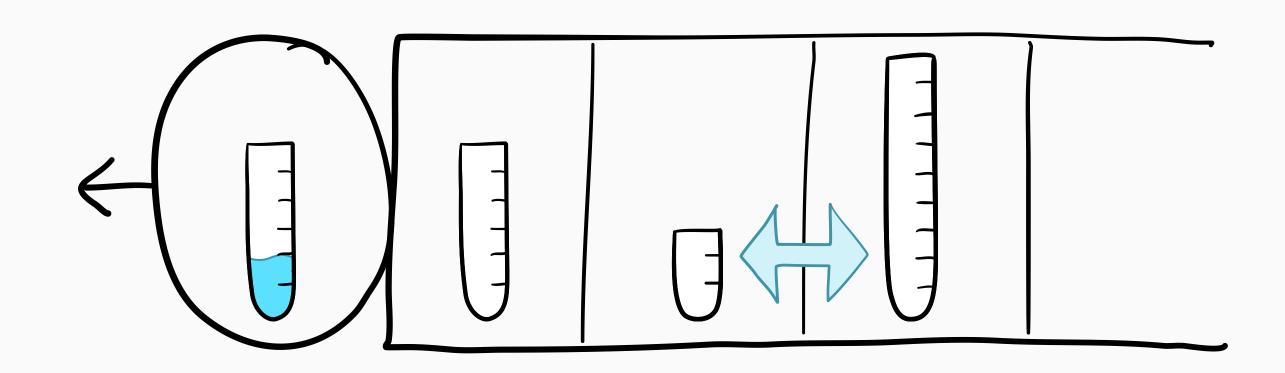


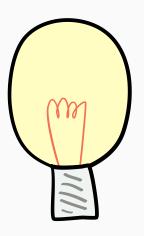
Nudge [Grosof et al., 2021]



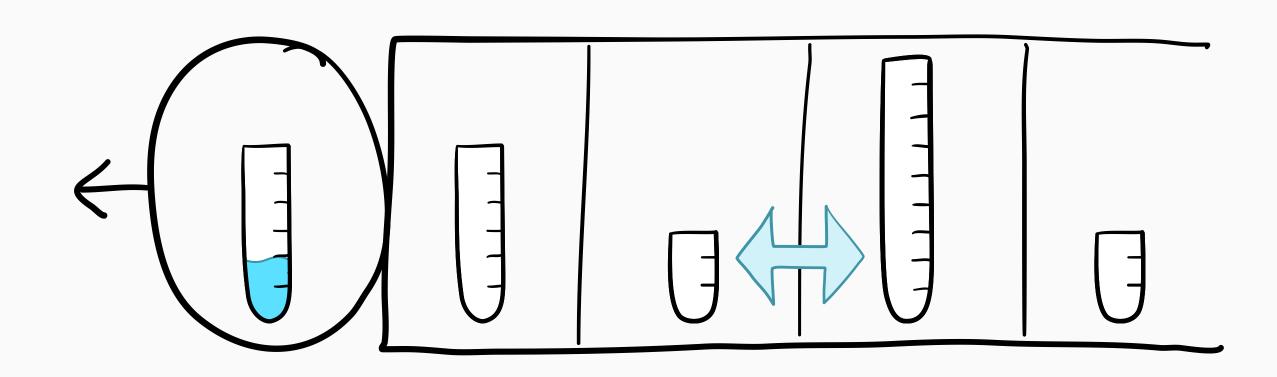


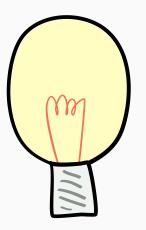
Nudge [Grosof et al., 2021]



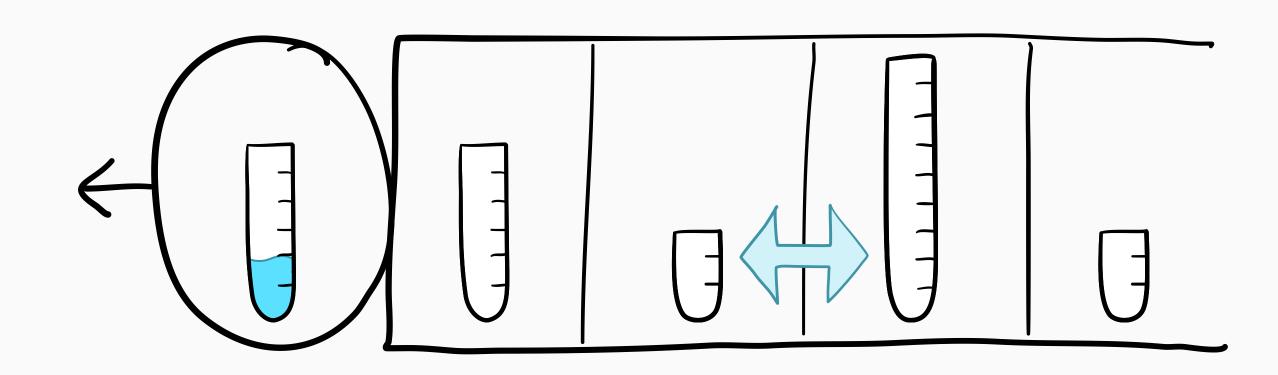


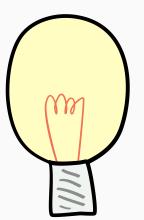
Nudge [Grosof et al., 2021]





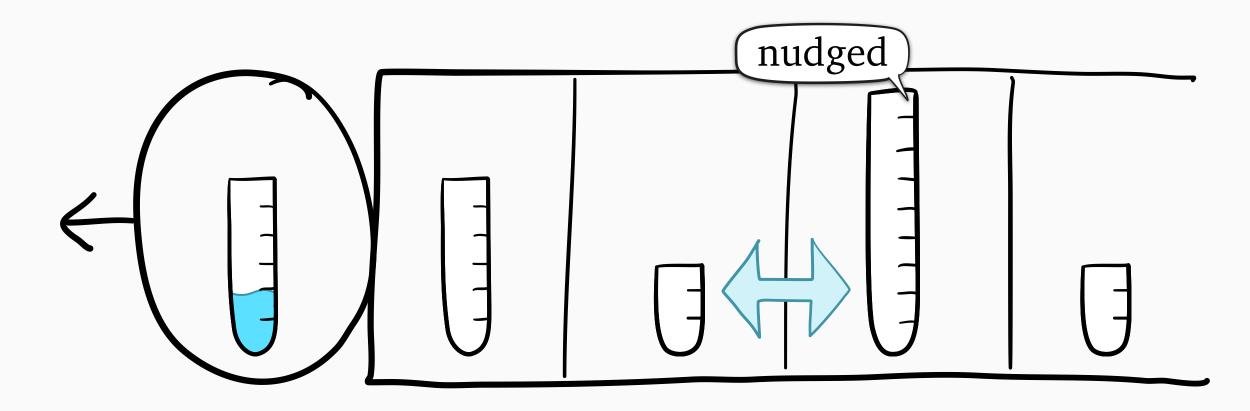
Nudge [Grosof et al., 2021]

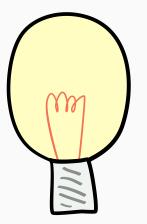




#### Nudge [Grosof et al., 2021]

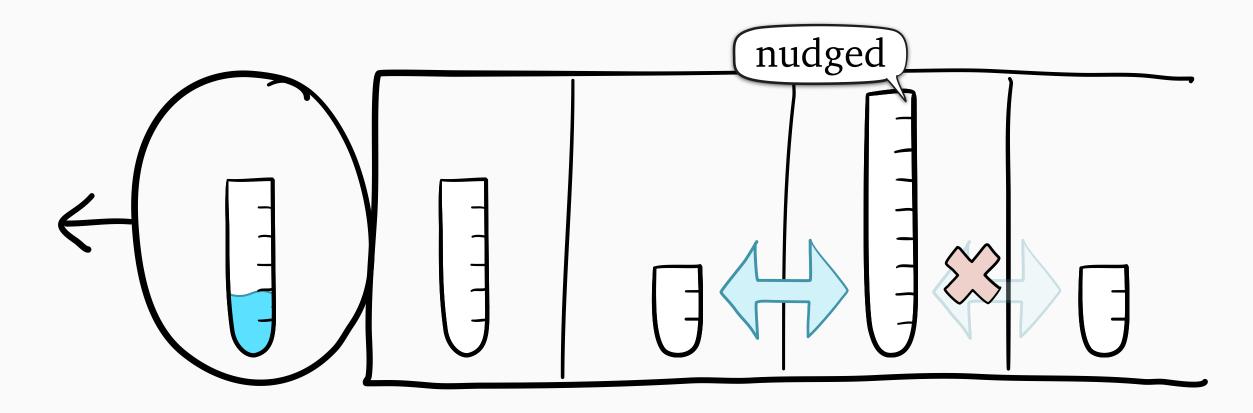
- small job can pass one large job
- large job can't be passed twice

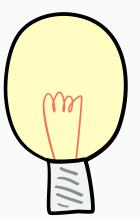




#### Nudge [Grosof et al., 2021]

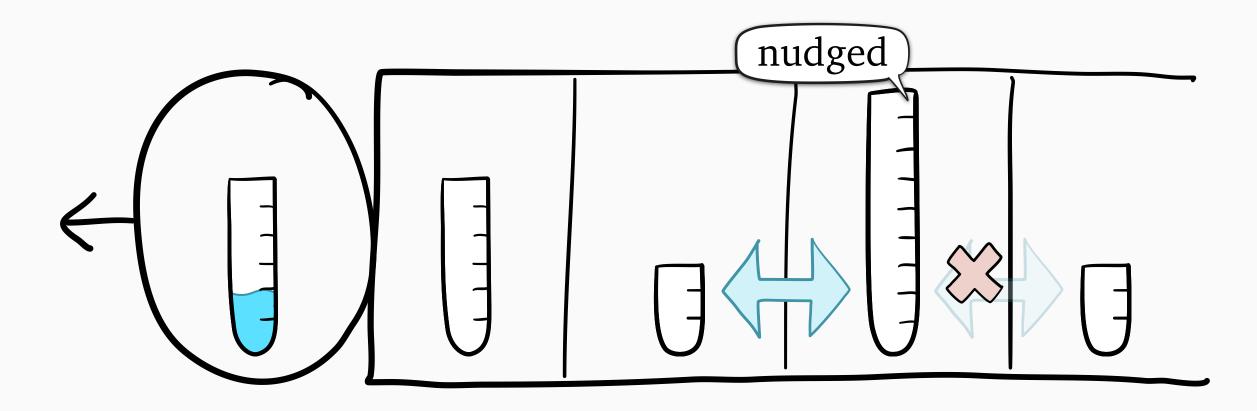
- small job can pass one large job
- large job can't be passed twice

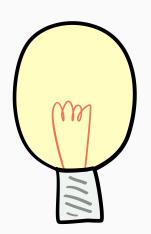




#### Nudge [Grosof et al., 2021]

- small job can pass one large job
- large job can't be passed twice



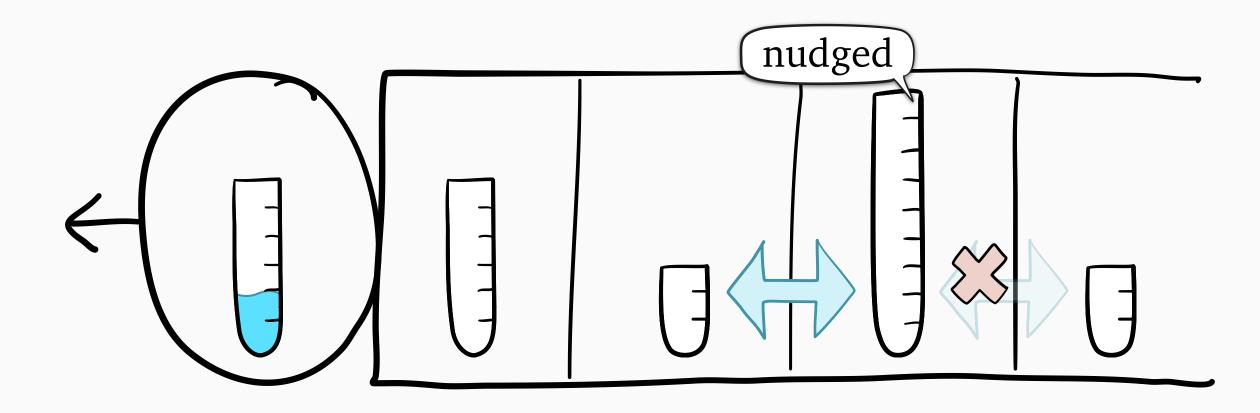


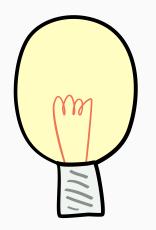
#### Nudge [Grosof et al., 2021]

- small job can pass one large job
- large job can't be passed twice

#### Theorem:

 $C_{\rm Nudge} < C_{\rm FCFS}$ 





#### Nudge [Grosof et al., 2021]

- small job can pass one large job
- large job can't be passed twice

#### Theorem:

 $C_{\rm Nudge} < C_{\rm FCFS}$ 

More complex variants get even lower C

[Van Houdt, 2022; Charlet & Van Houdt, 2024]

# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



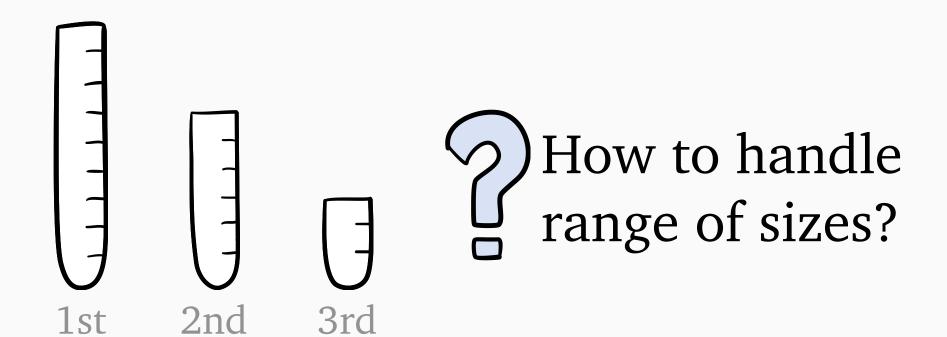
Why is achieving strong tail optimality hard?

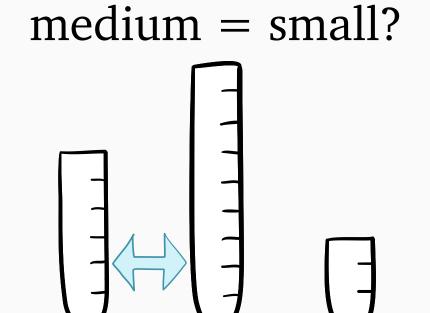


How does the **Boost** policy family work?

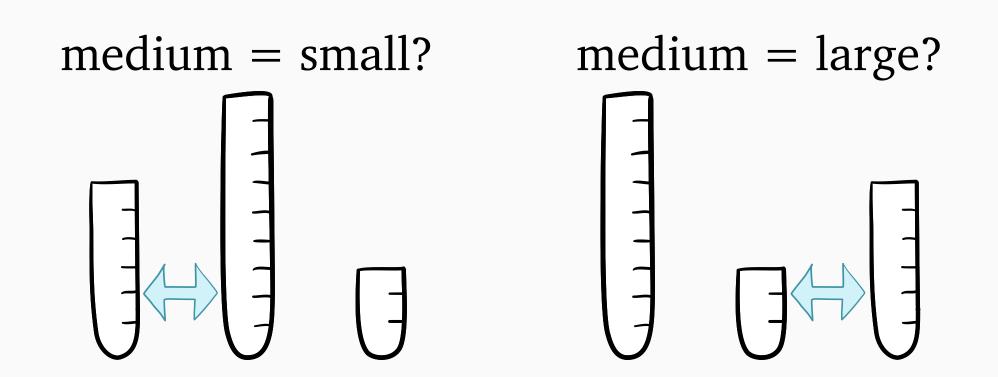


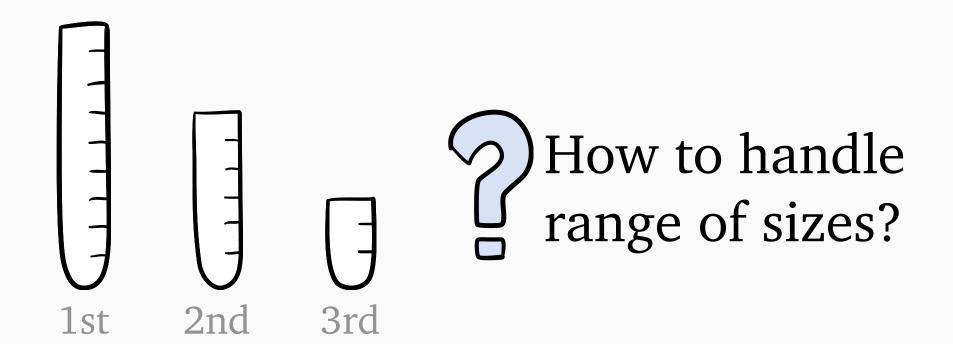


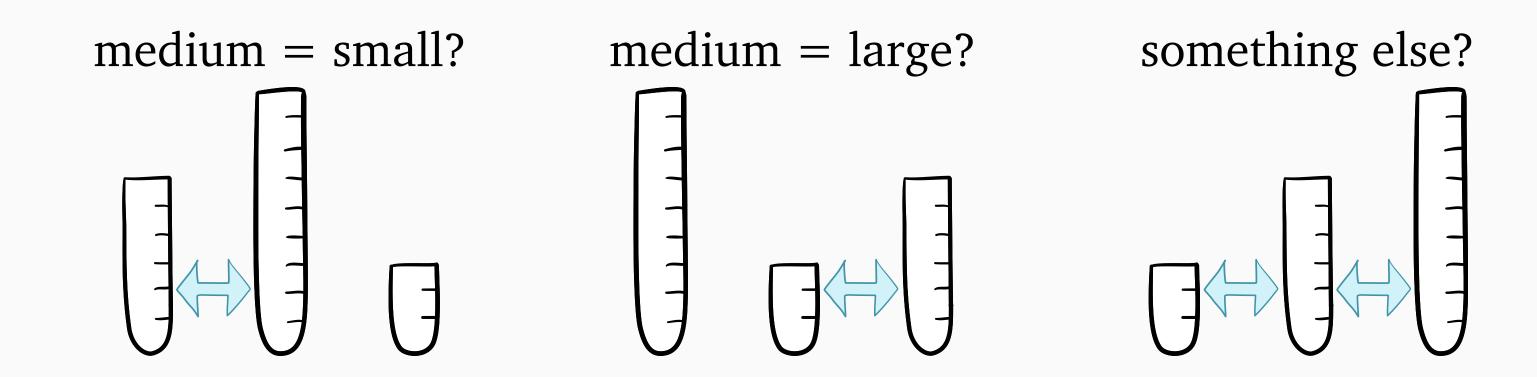


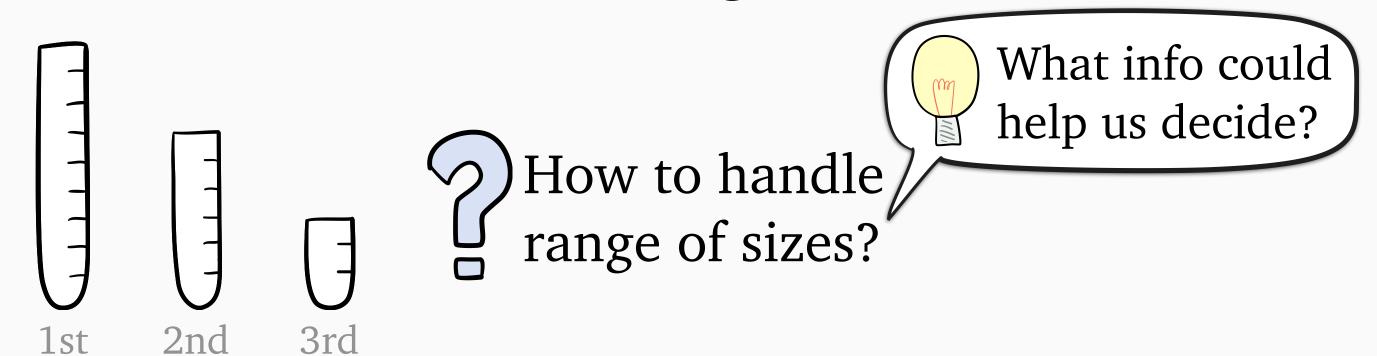


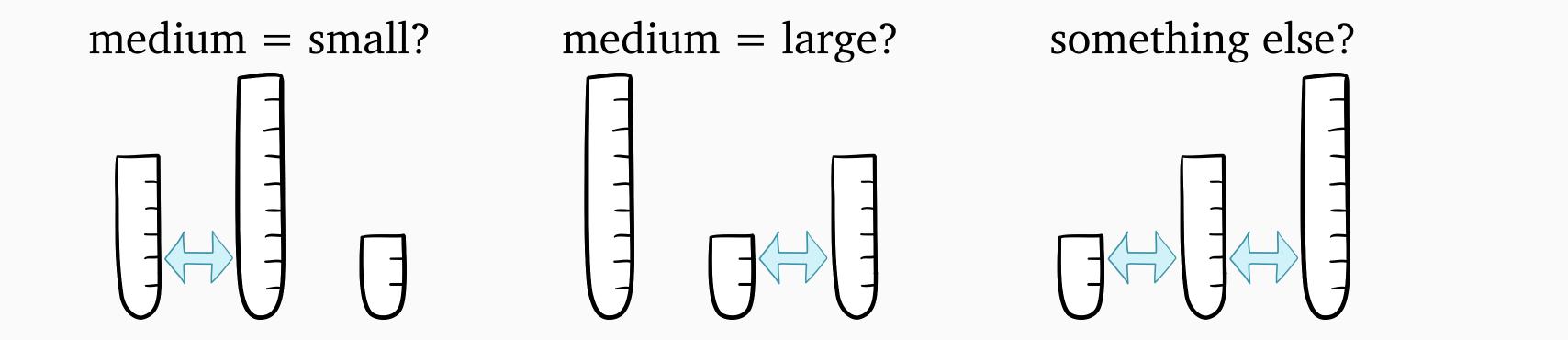


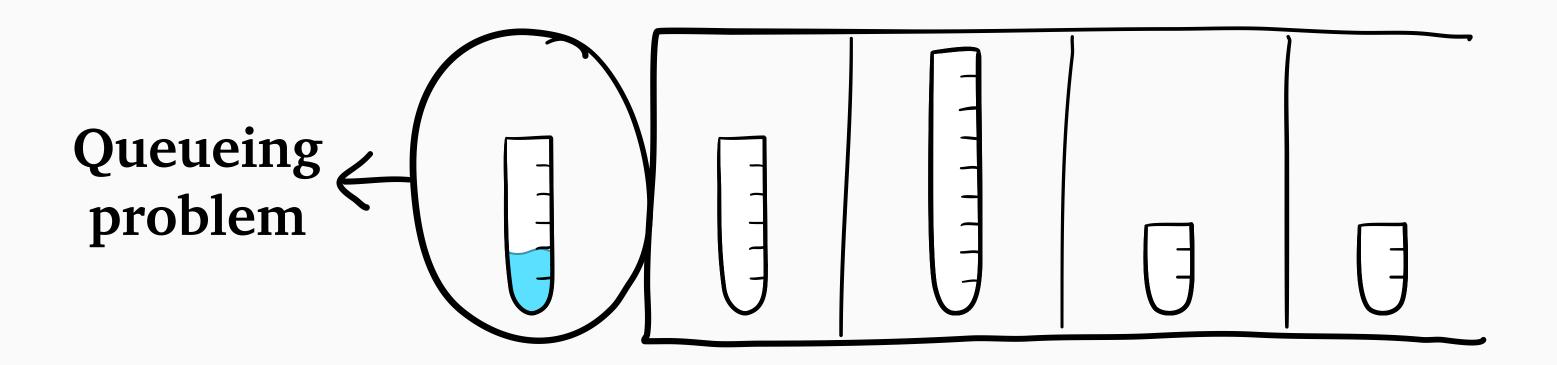


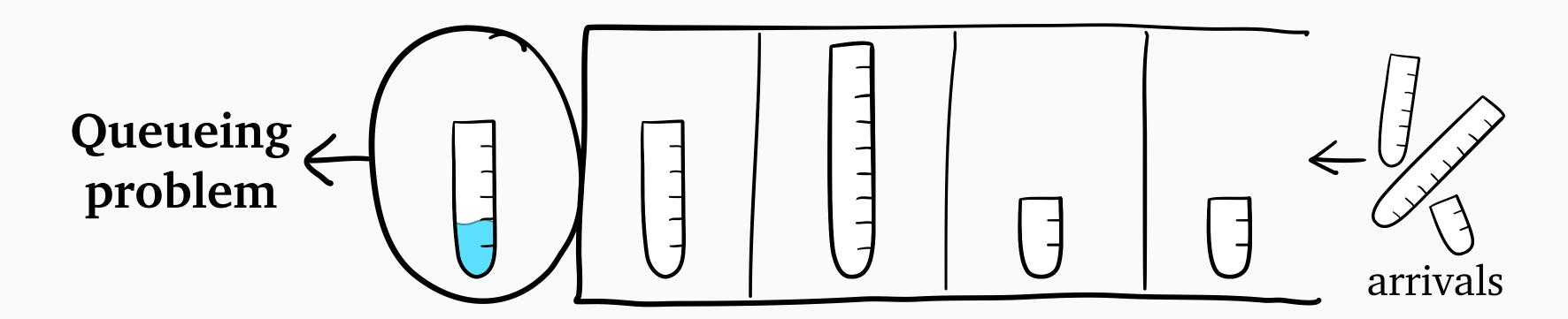


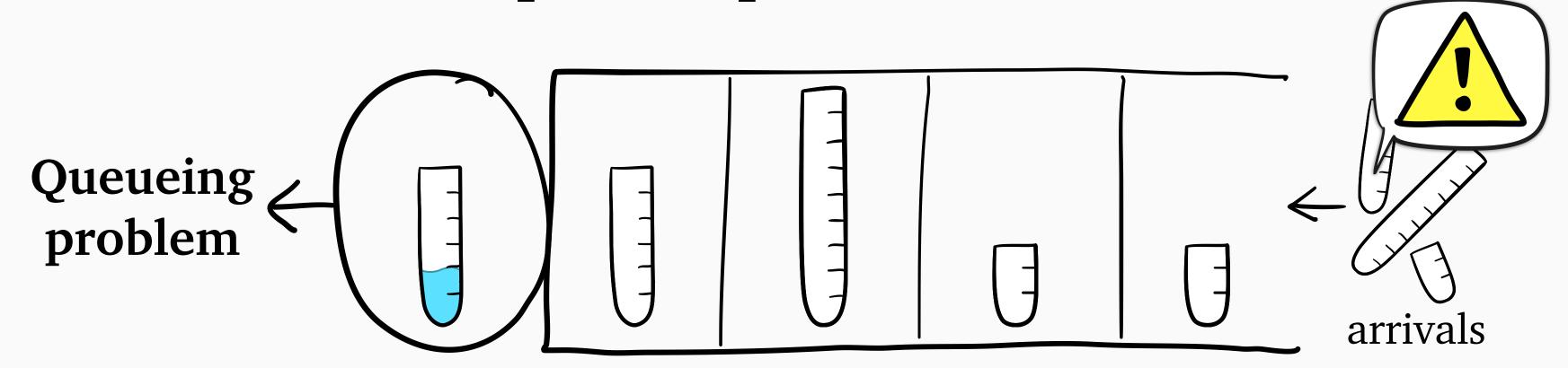


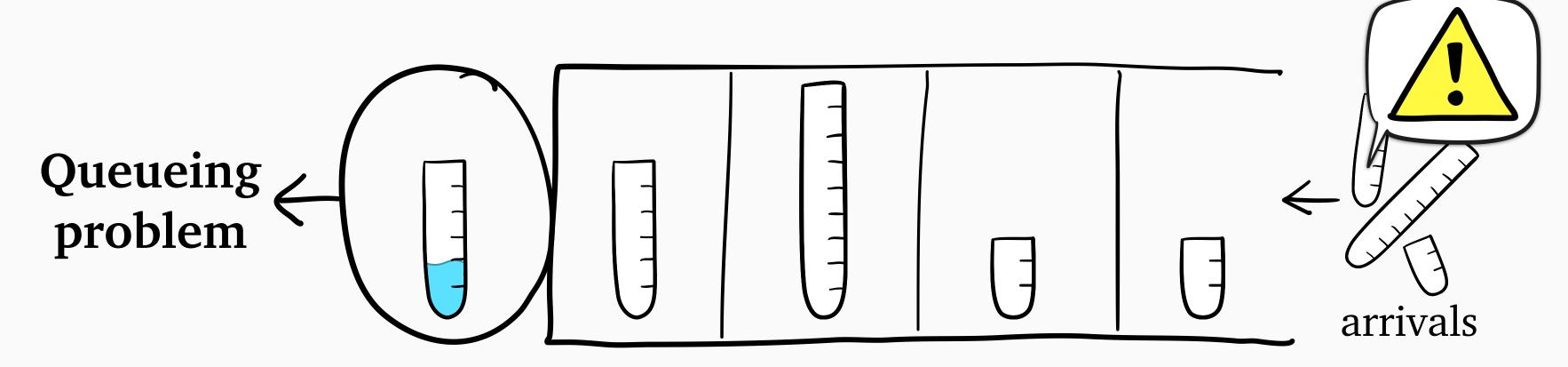


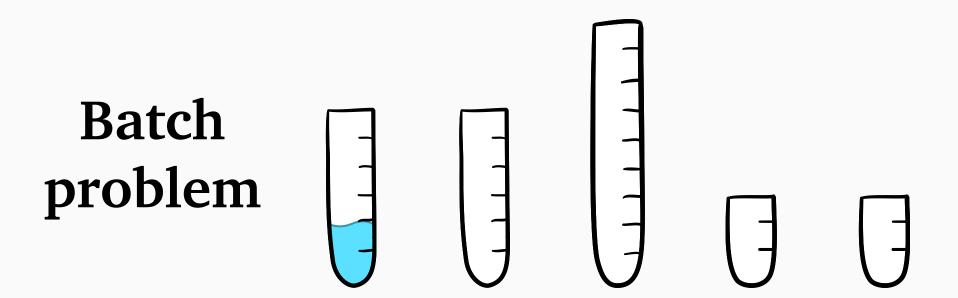


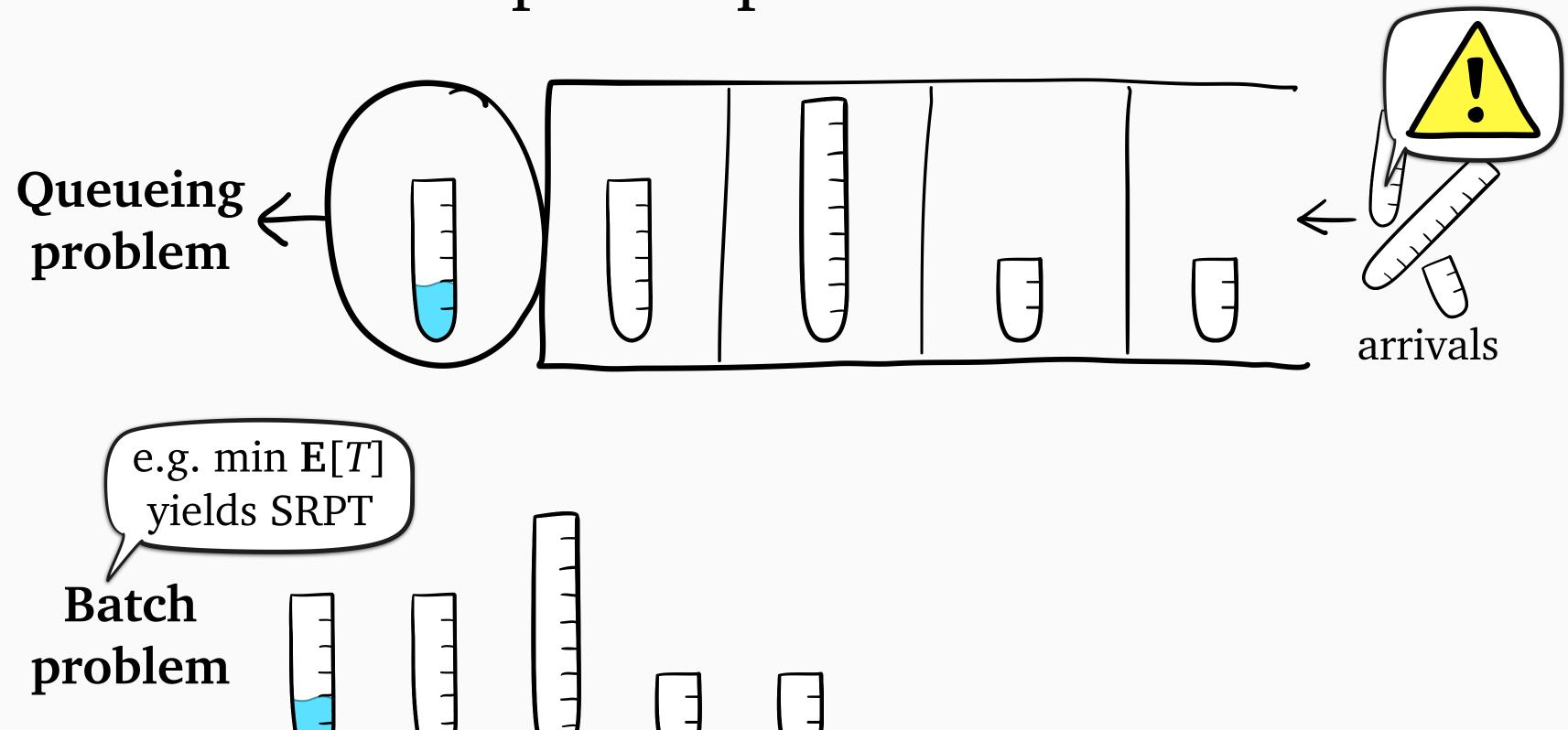


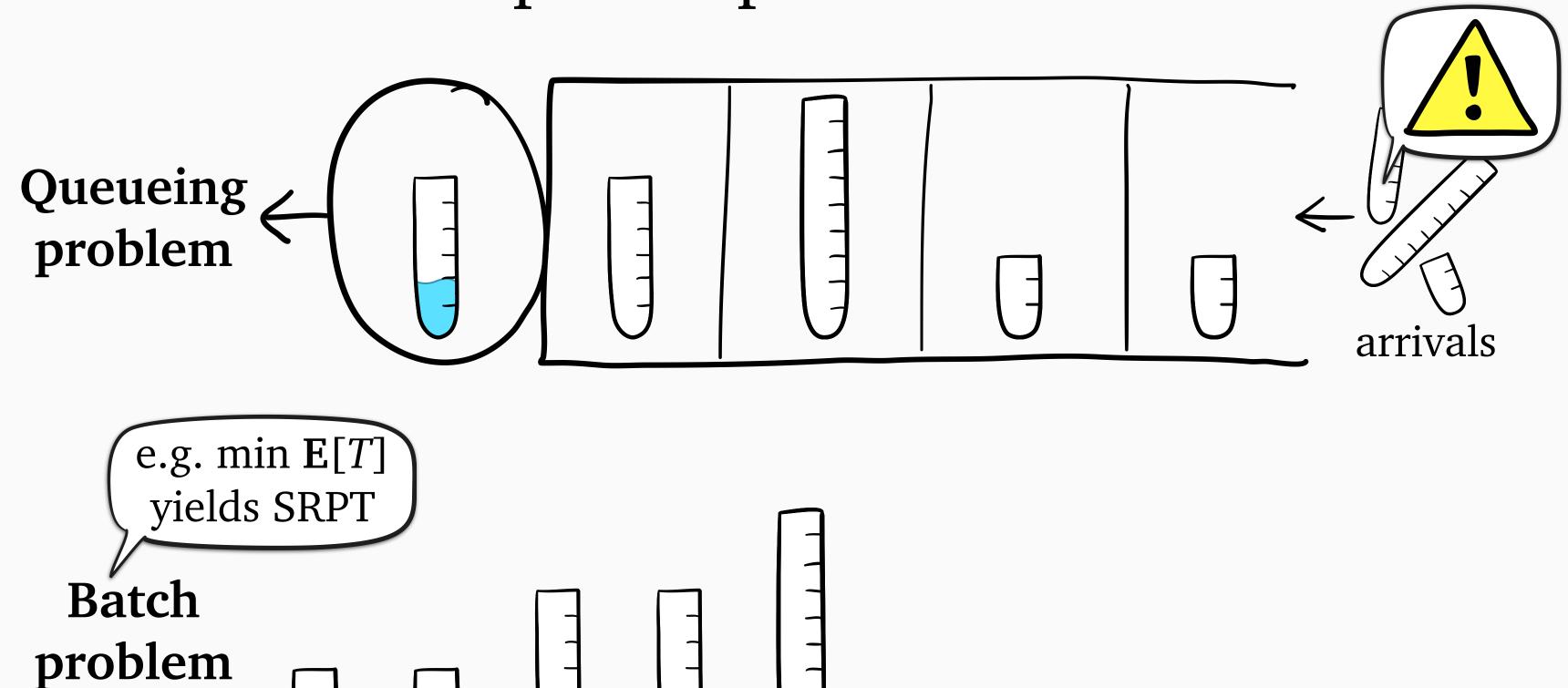


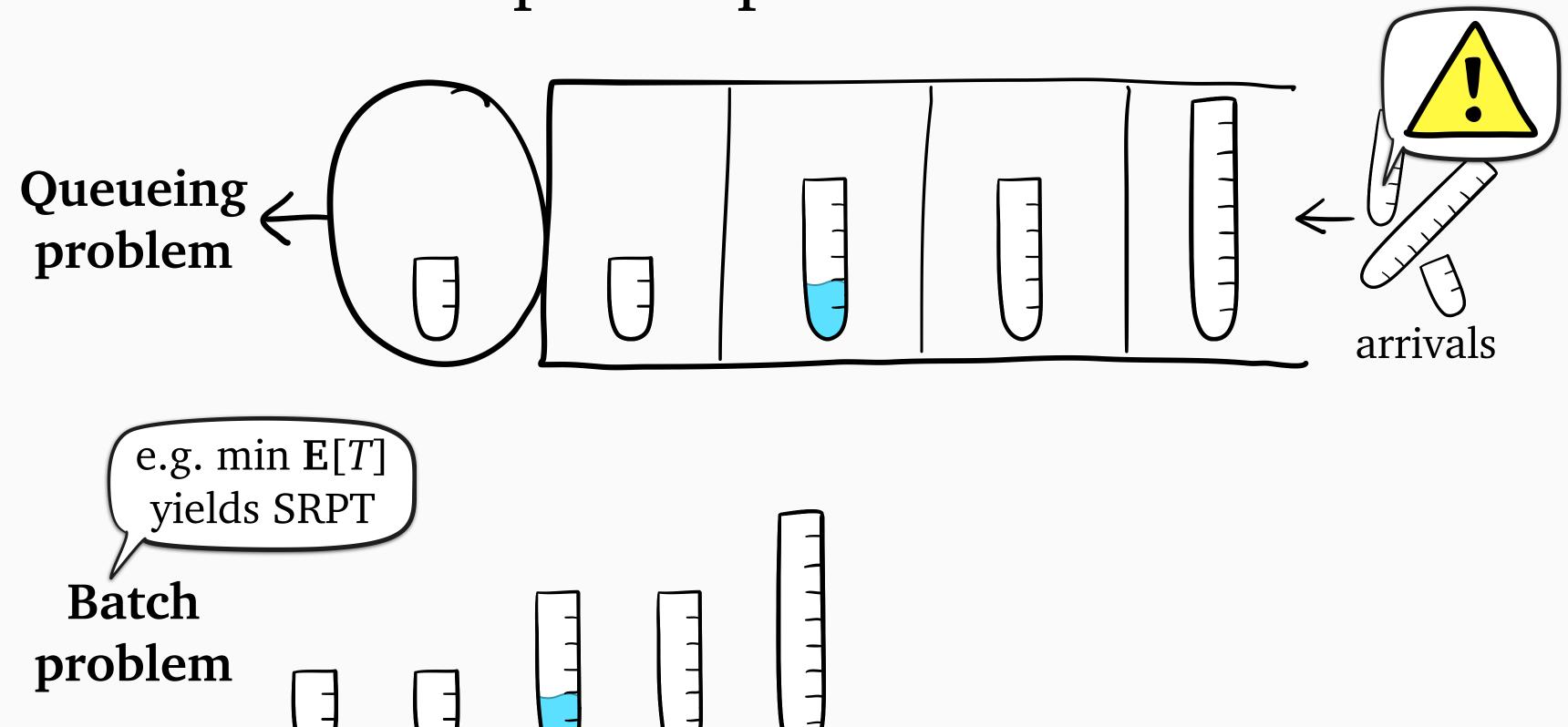


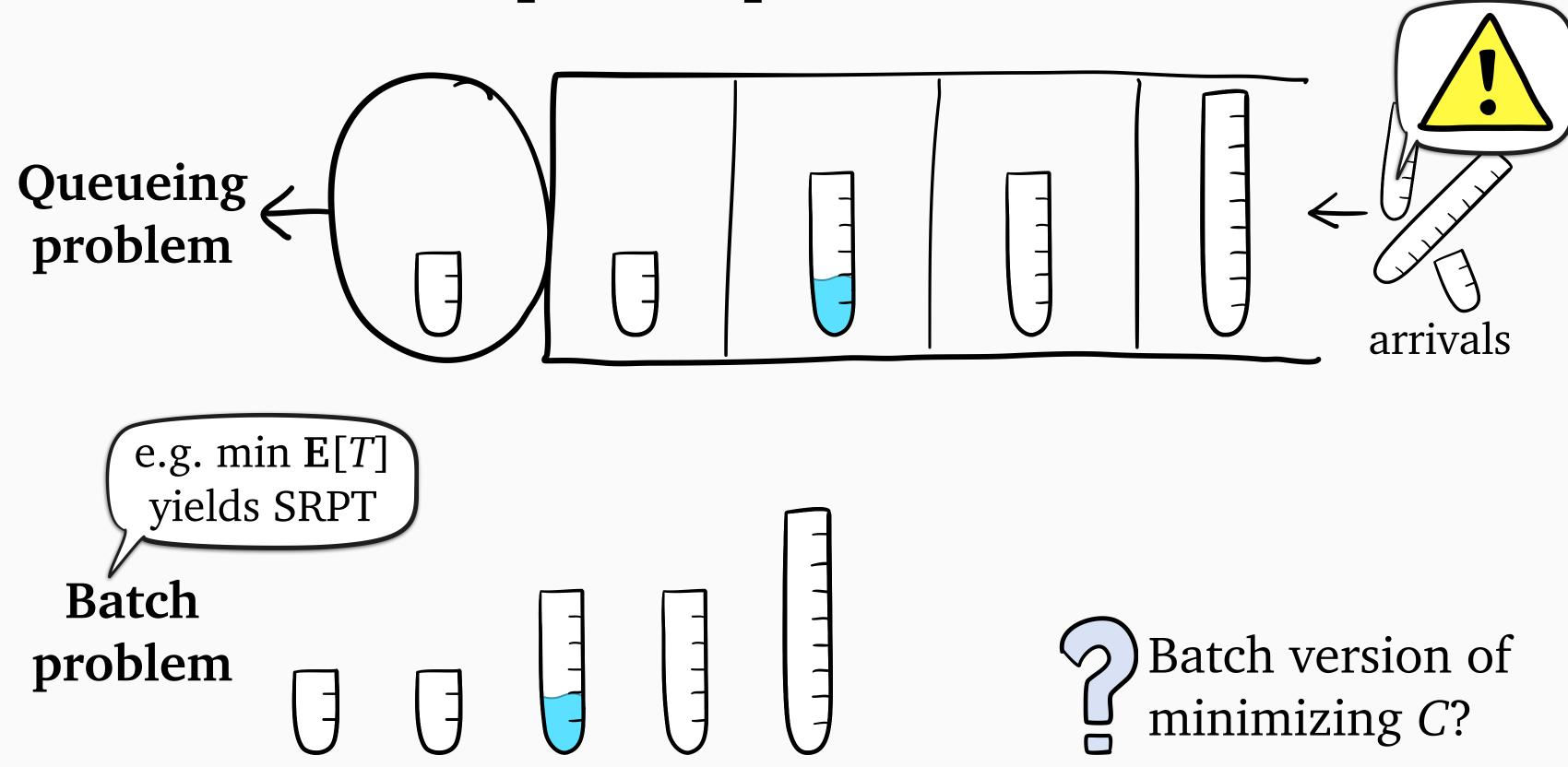


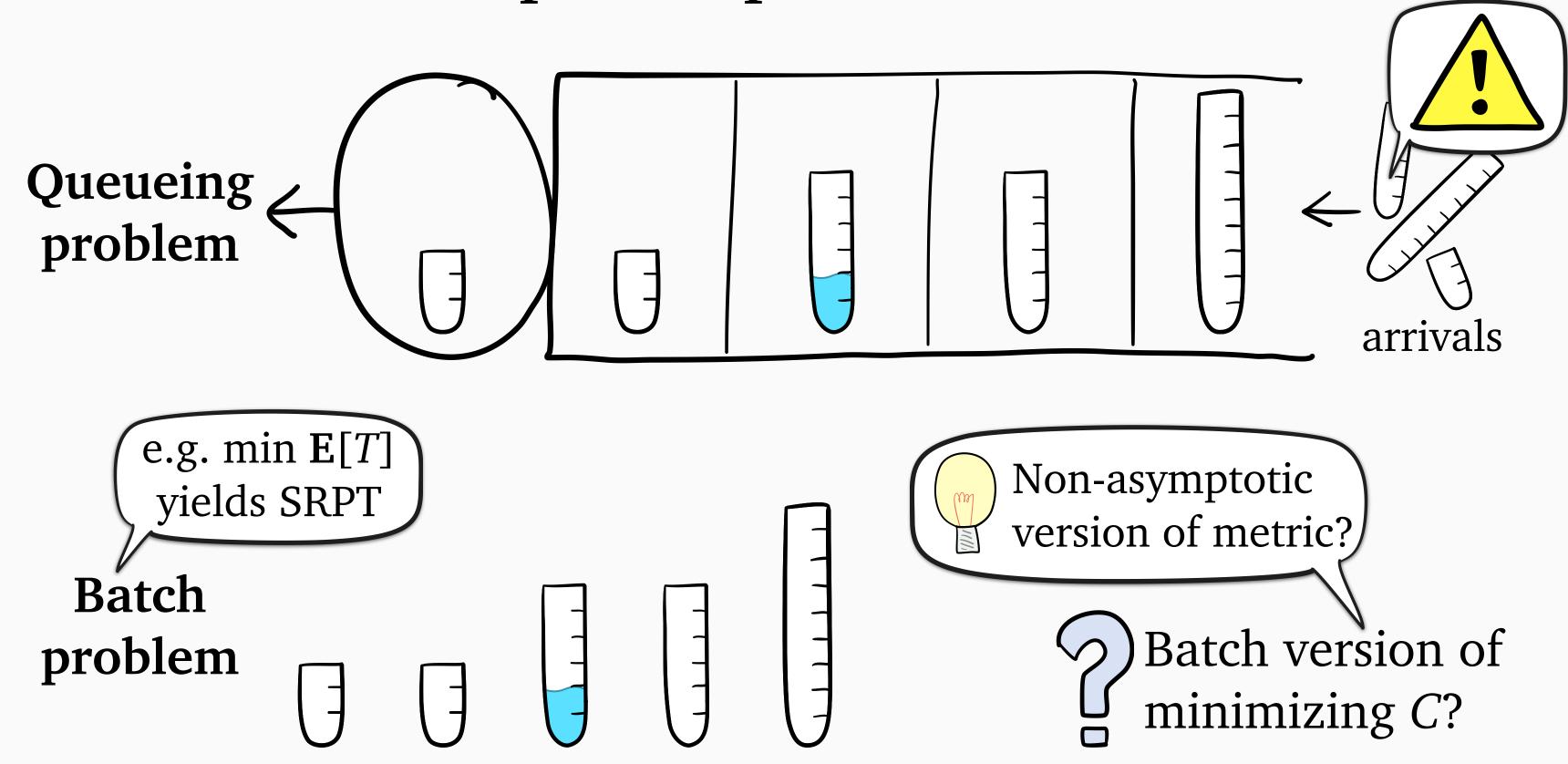












# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



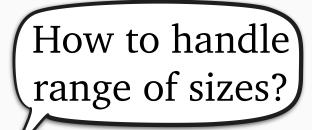
How does the **Boost** policy family work?



# Boost



Why did it take so long to beat FCFS?





Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?

Batch version of minimizing *C*?



# Boost



Why did it take so long to beat FCFS?





Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?

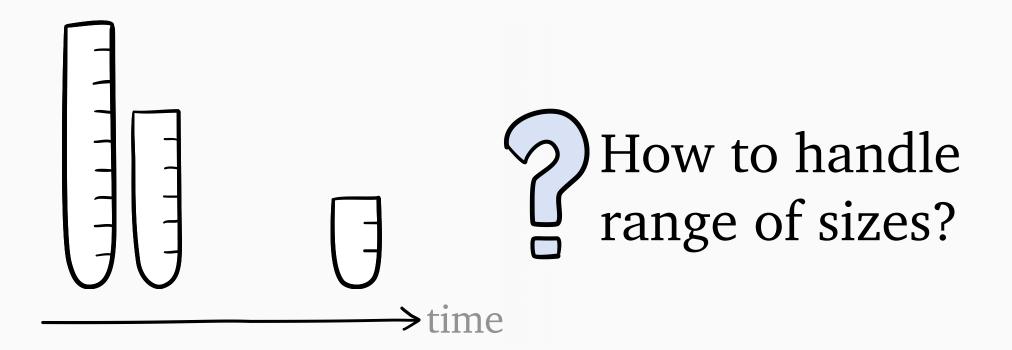
Batch version of minimizing *C*?

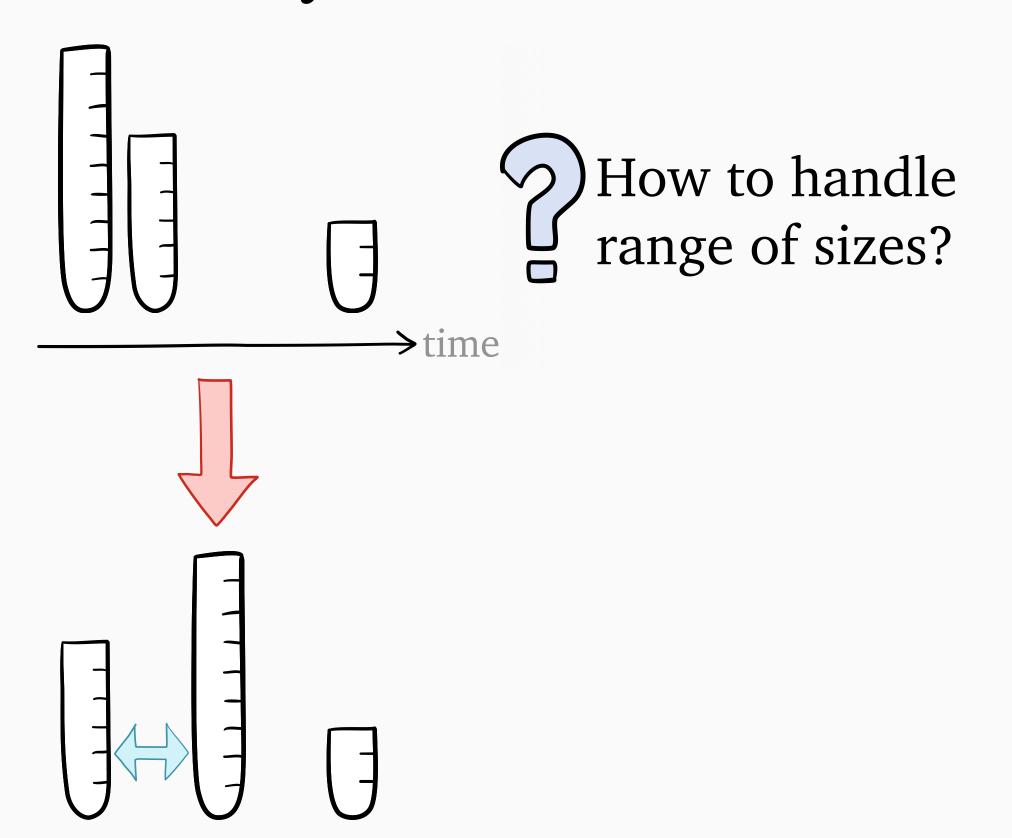


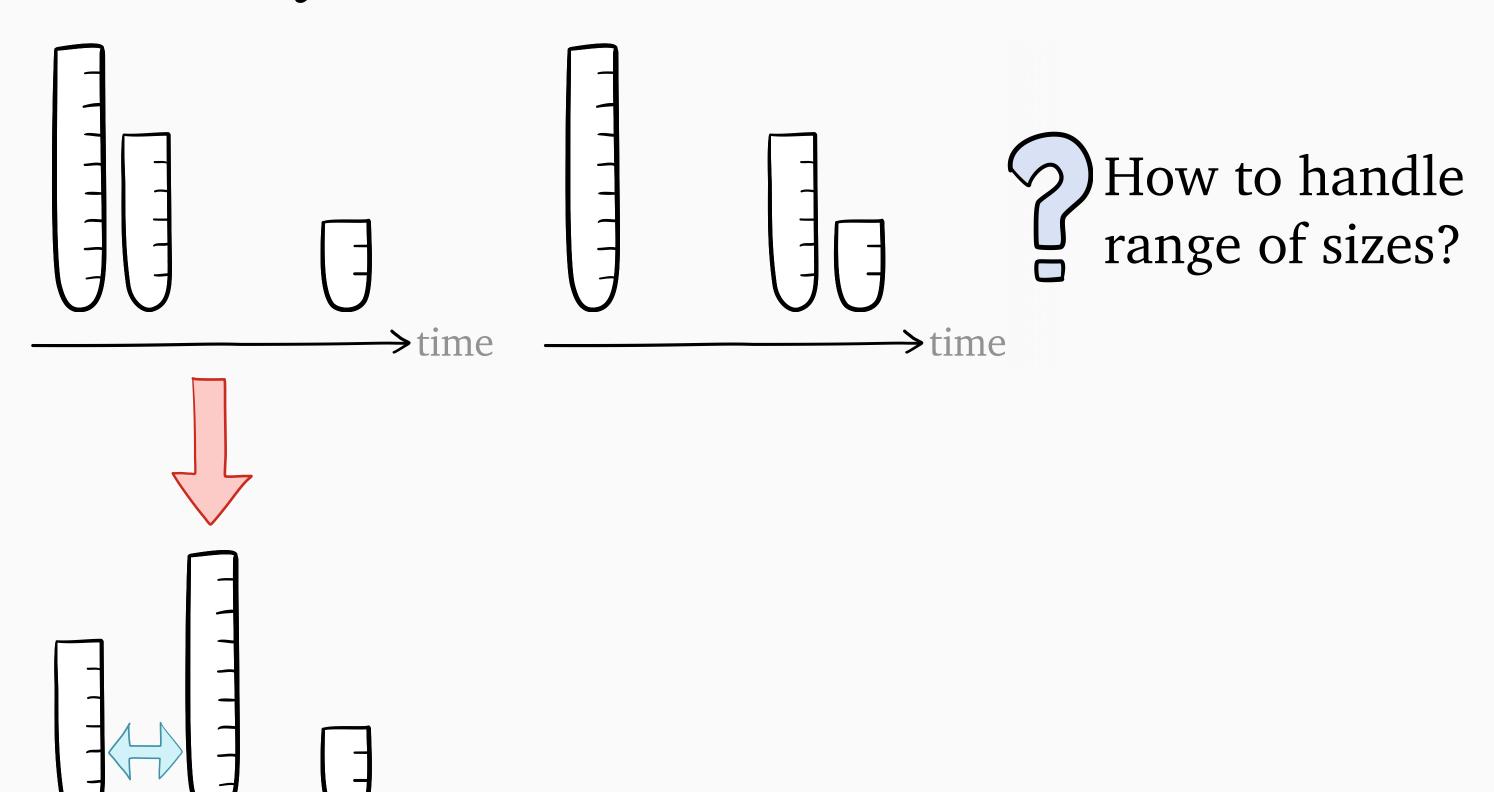
### Key information:

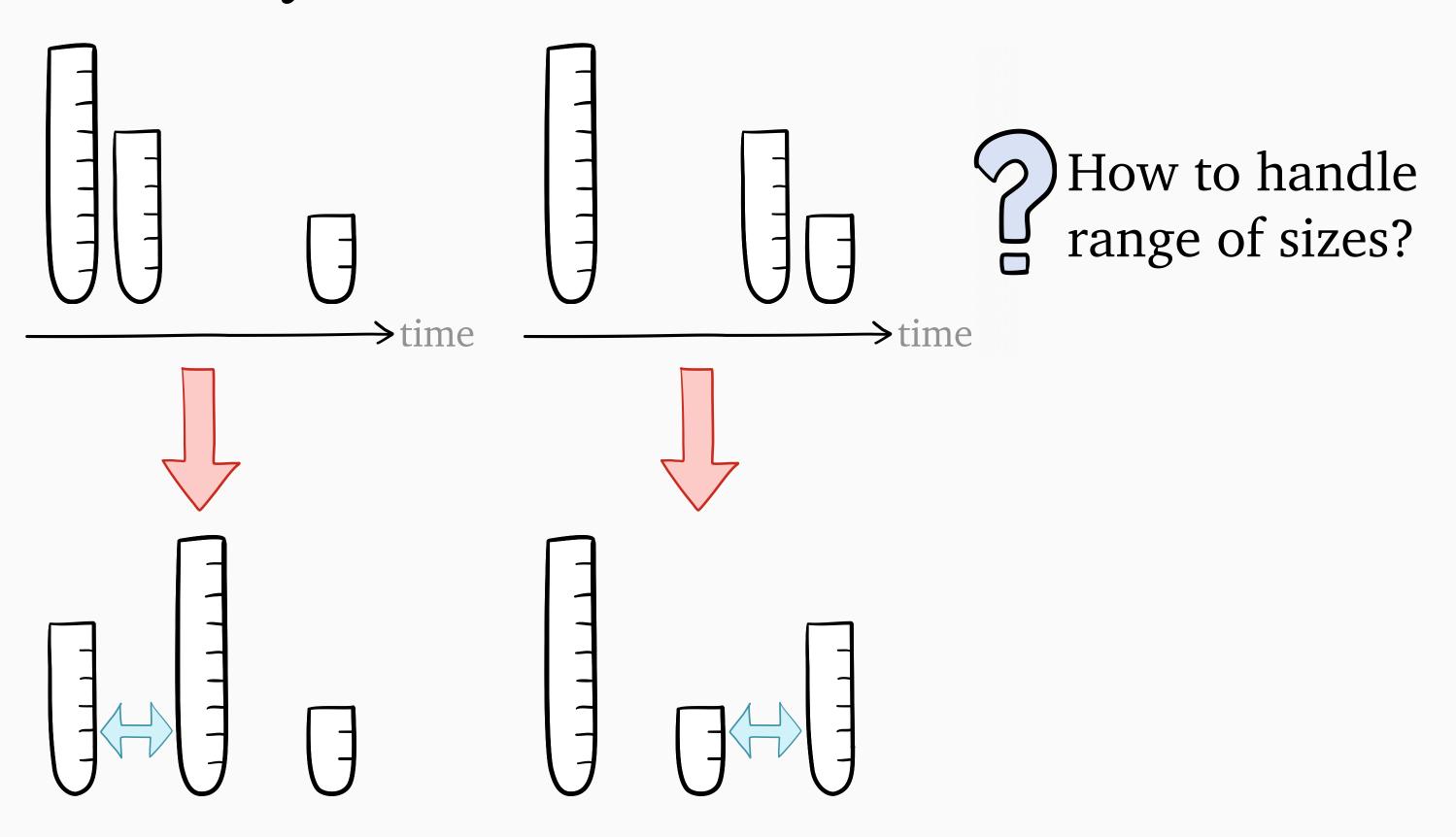


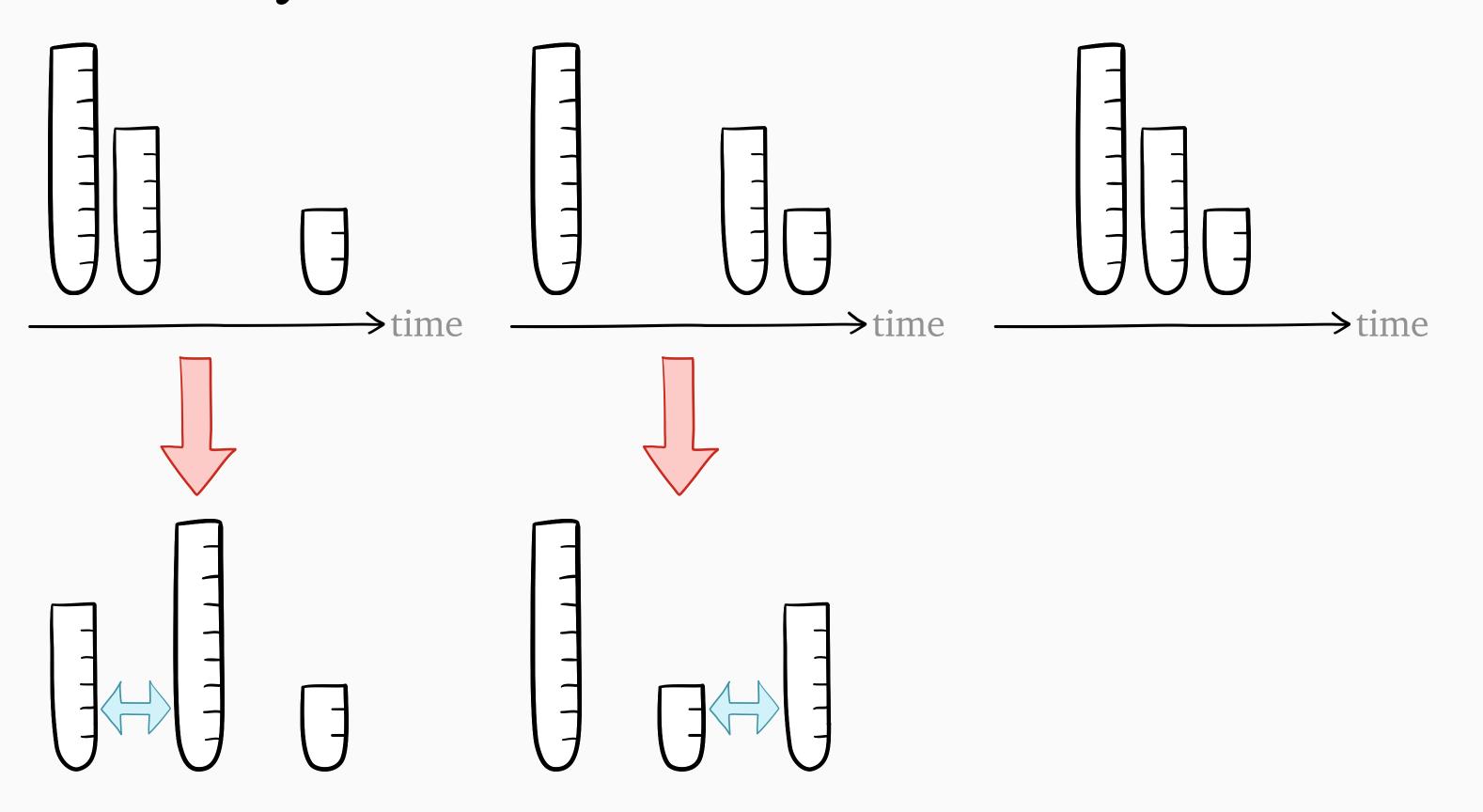


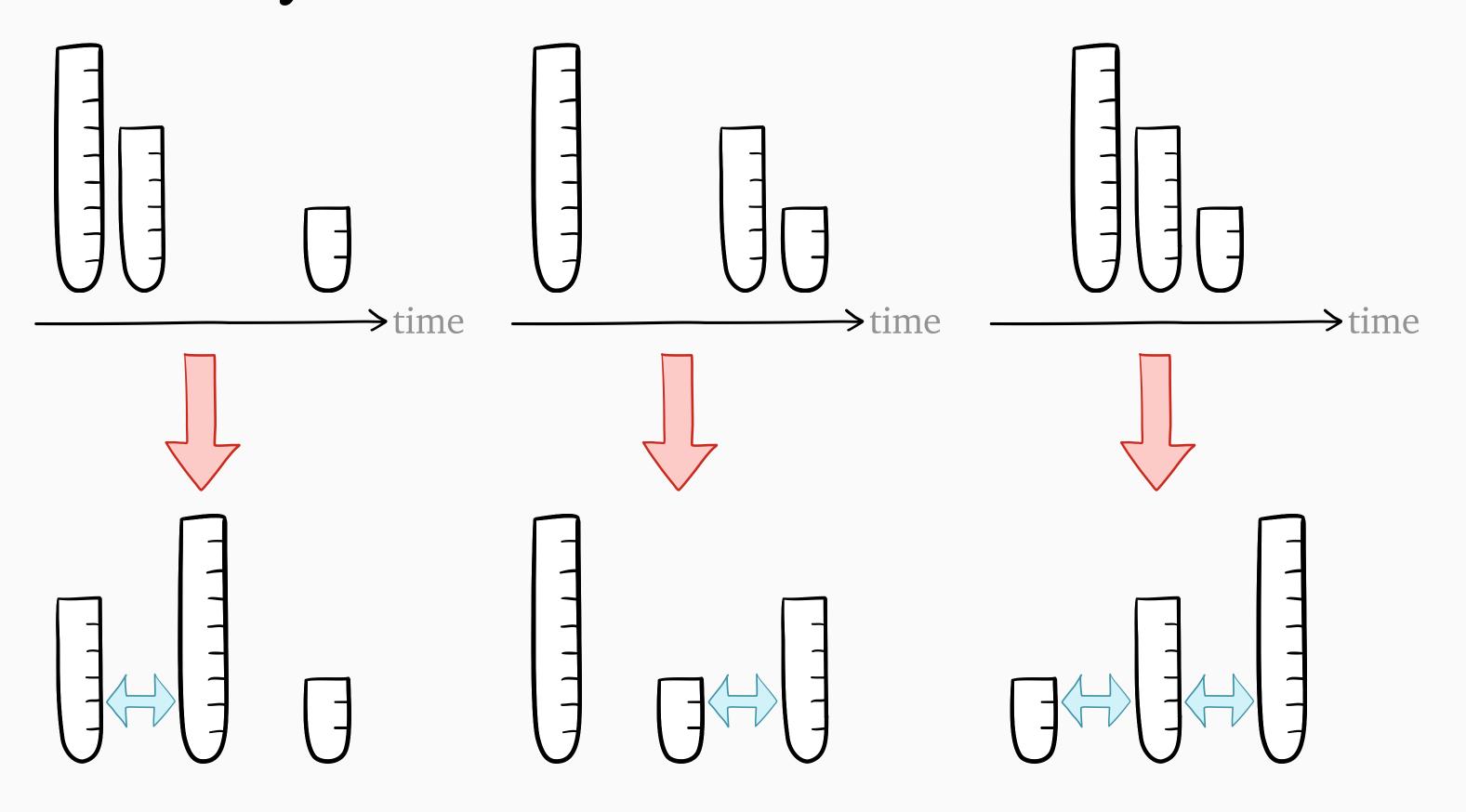


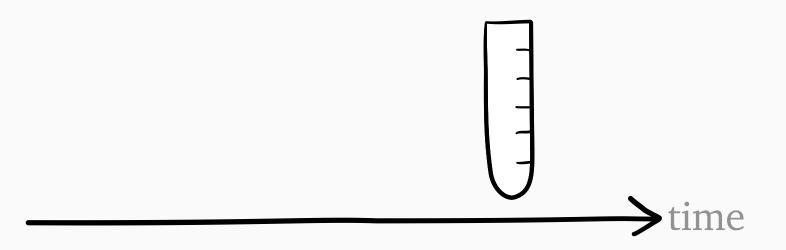


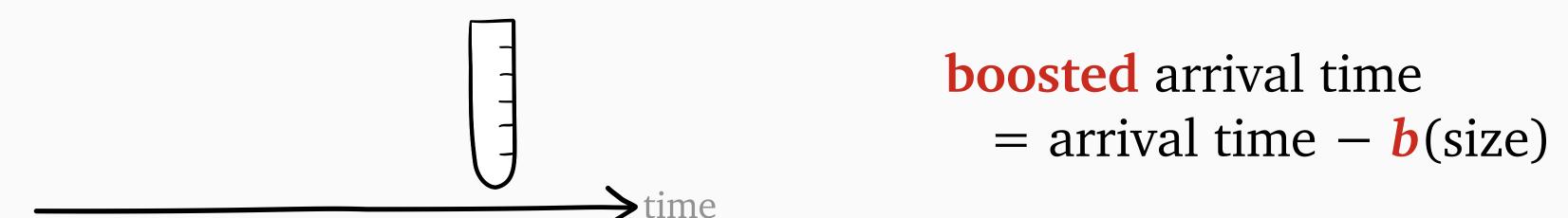


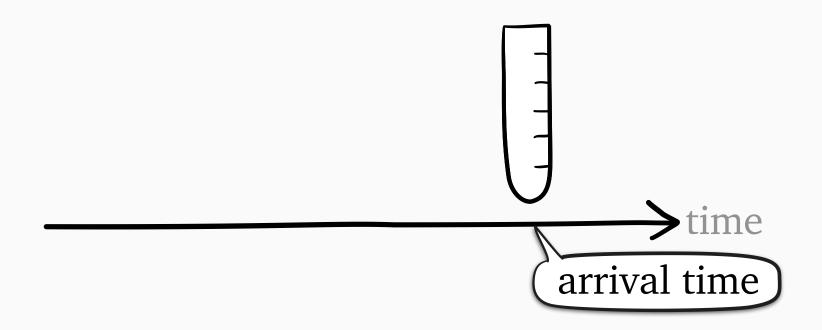




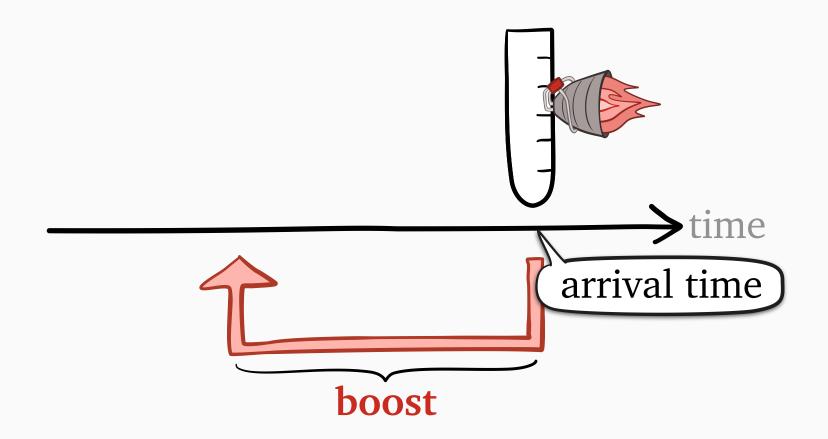




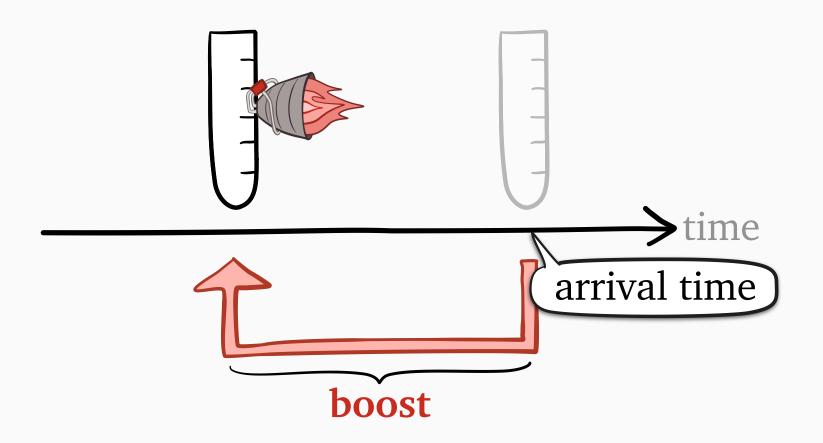




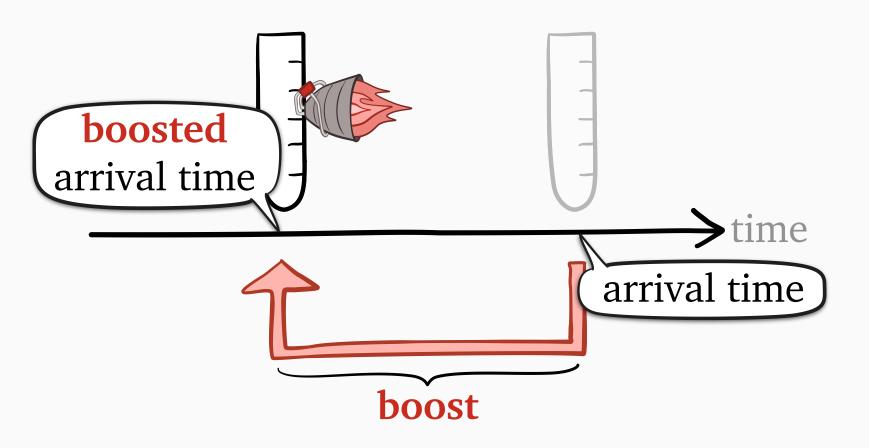
**boosted** arrival time



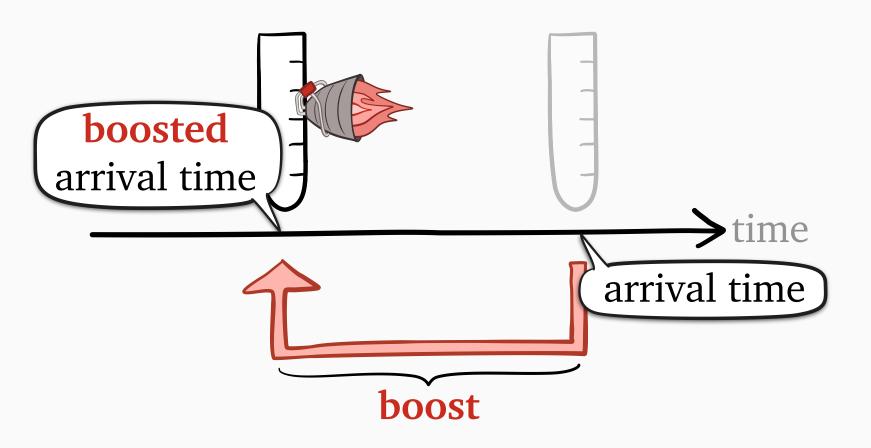
**boosted** arrival time



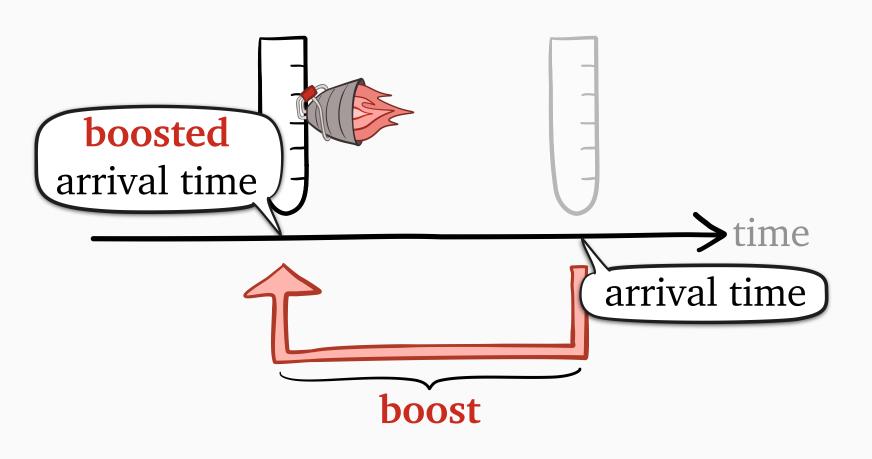
**boosted** arrival time



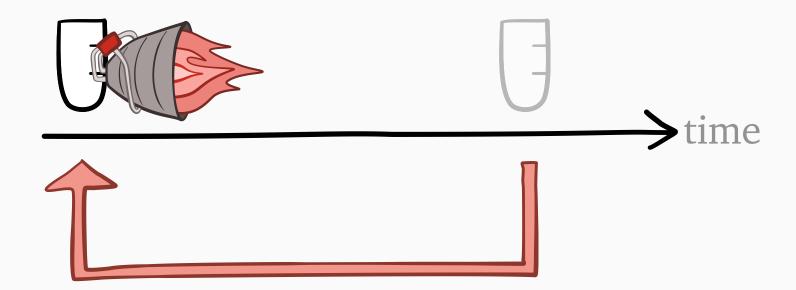
#### **boosted** arrival time

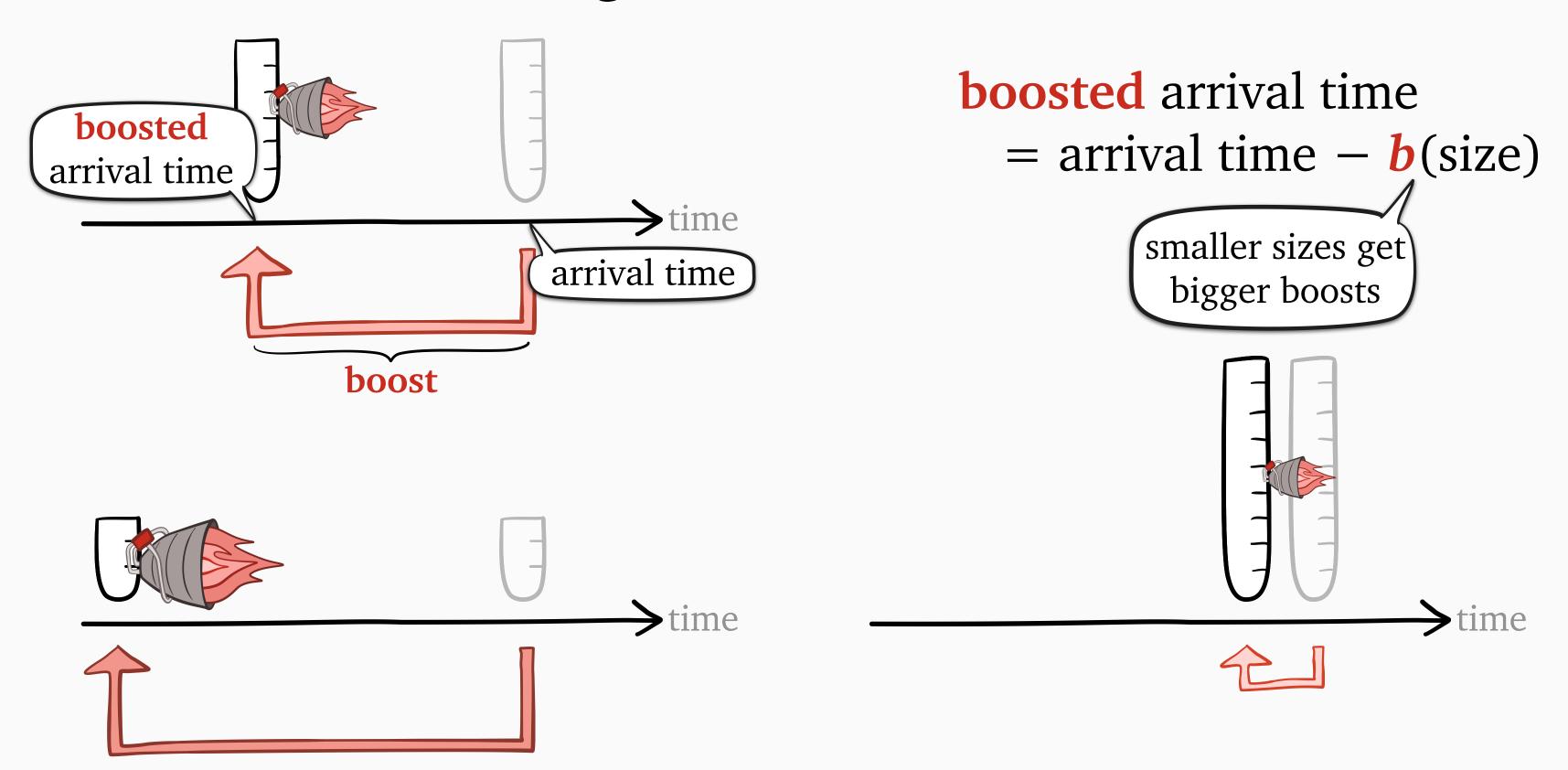


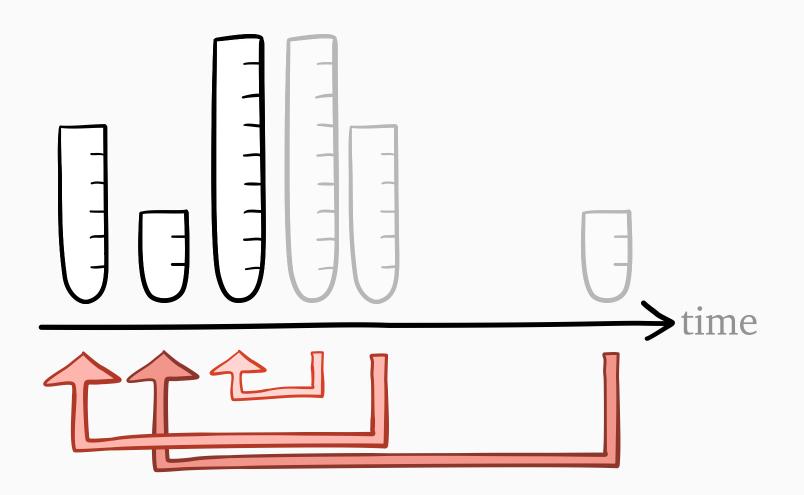








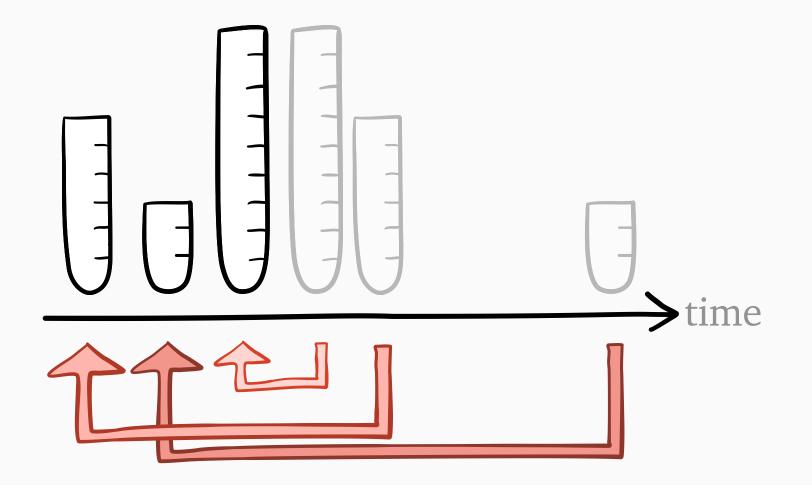




**boosted** arrival time



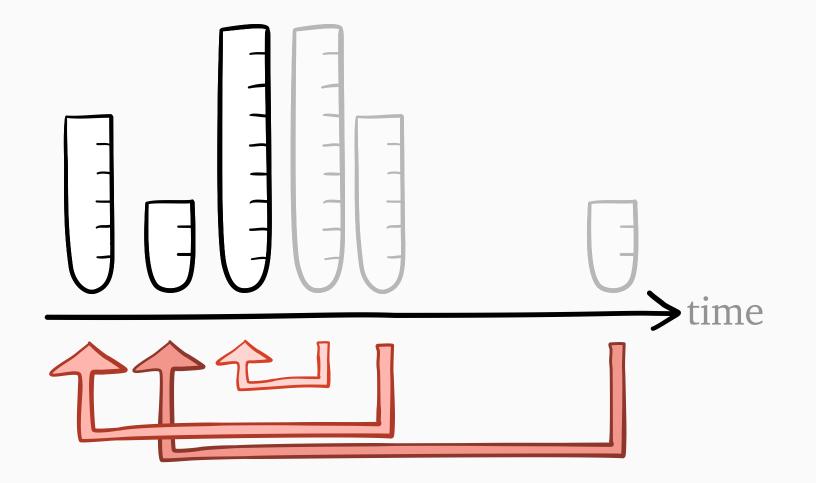
Scheduling rule: always serve job of minimum boosted arrival time

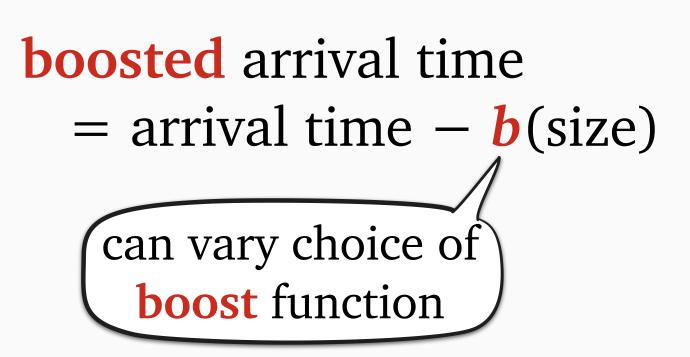


**boosted** arrival time



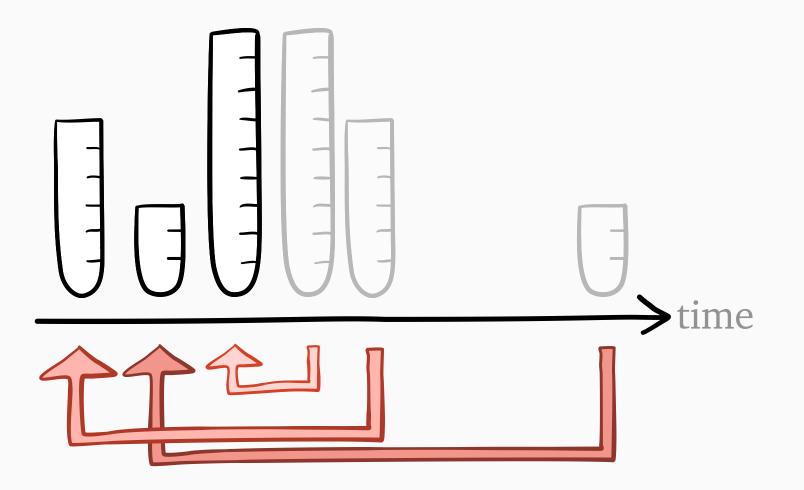
Scheduling rule: always serve job of minimum boosted arrival time





can be preemptive or nonpreemptive

Scheduling rule: always serve job of minimum boosted arrival time



boosted arrival time
= arrival time - b(size)

can vary choice of
boost function

# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



# Boost



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



How do we achieve strong tail optimality?

What's the right boost function?

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

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



$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem



$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$

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

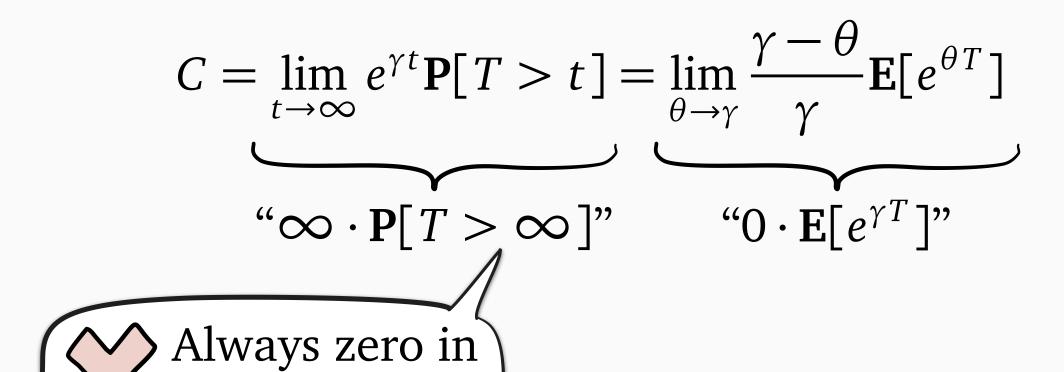


$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$

$$"\infty \cdot \mathbf{P}[T > \infty]"$$
Always zero in

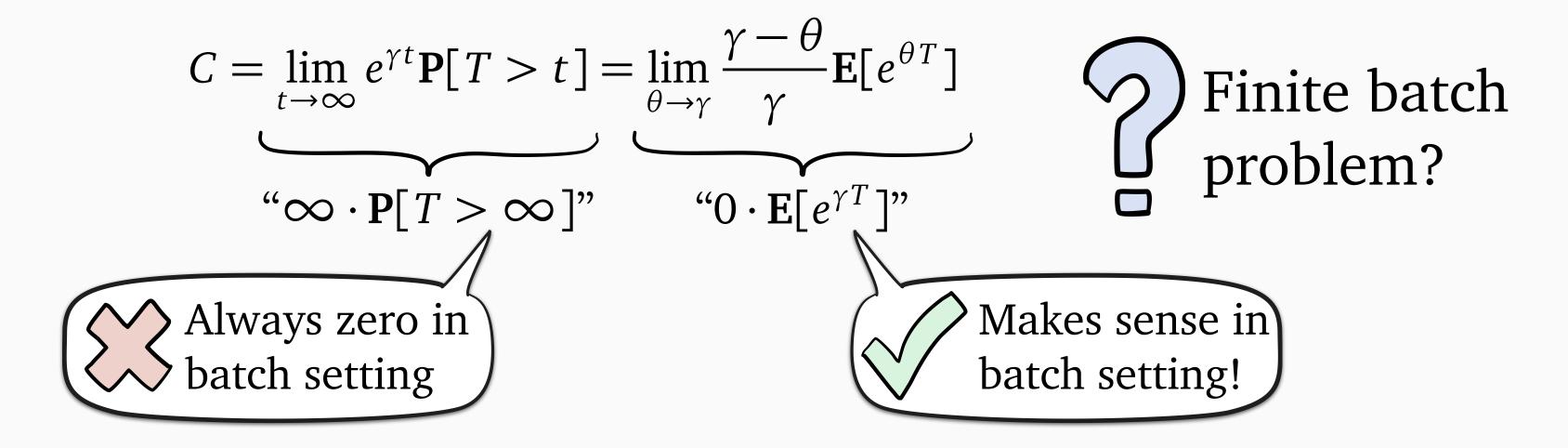
batch setting





batch setting





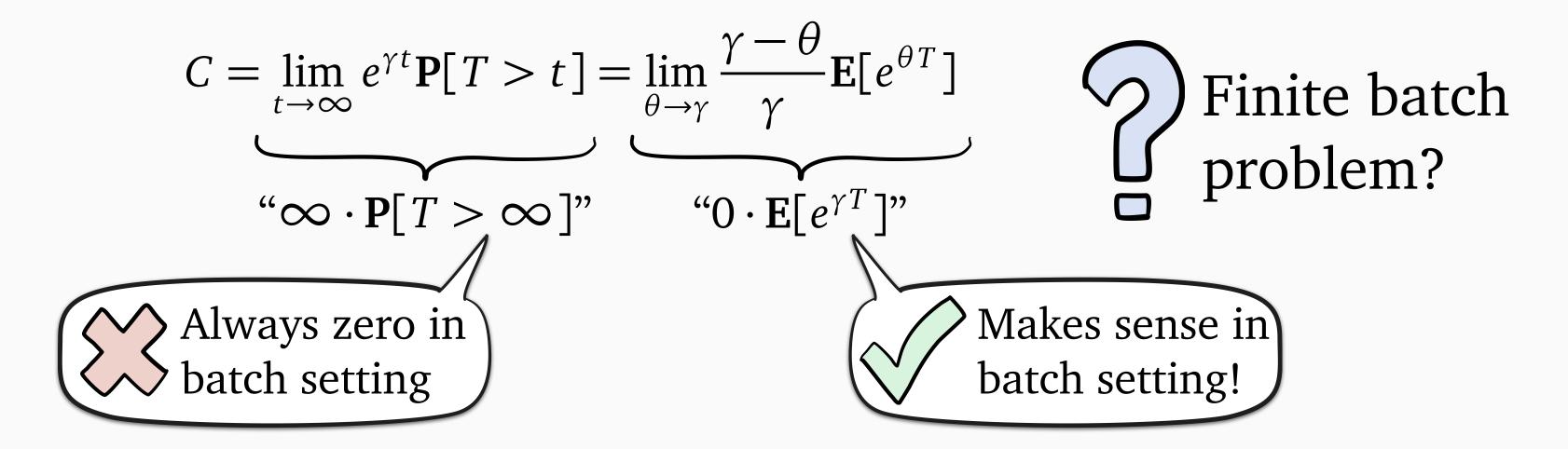
$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$

$$\text{Finite batch problem?}$$

$$(0) \times \mathbf{E}[e^{\gamma T}]$$

$$\text{Makes sense in batch setting!}$$

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



Batch problem: minimize

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



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

## Transforming the problem

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$

$$\text{Finite batch problem?}$$

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

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

Always zero in batch setting

Makes sense in batch setting!

almost classic problem

Batch problem: minimize

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



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

Batch problem: minimize

$$t_i = d_i - a_i$$

 $d_i$  = departure time of job i



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

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

Batch problem: minimize

$$t_i = d_i - a_i$$

 $d_i$  = departure time of job i



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

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

Classic metric: mean weighted discounted departure time

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

Batch problem: minimize

$$t_i = d_i - a_i$$

 $d_i$  = departure time of job i



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

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



Classic metric: mean weighted  $\frac{1}{n} \sum_{i=1}^{n} w_i e^{-\theta d_i}$  discounted departure time

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

Batch problem: minimize

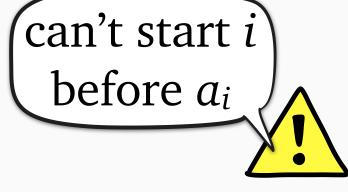
$$t_i = d_i - a_i$$

 $a_i$  = arrival time of job i

 $d_i$  = departure time of job i



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



Classic metric: mean weighted discounted departure time

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

Batch problem: minimize

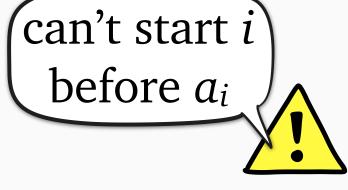
$$t_i = d_i - a_i$$

 $a_i$  = arrival time of job i

 $d_i$  = departure time of job i



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



Classic metric: mean weighted discounted departure time

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

Batch problem: minimize

$$t_i = d_i - a_i$$

 $a_i$  = arrival time of job i

 $d_i$  = departure time of job i



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

can't start i before  $a_i$ 

Classic metric: mean weighted discounted departure time

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

negative discount rate

Batch problem: minimize

$$t_i = d_i - a_i$$

 $a_i$  = arrival time of job i

 $d_i$  = departure time of job i



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

can't start i before  $a_i$ 

Classic metric: mean weighted discounted departure time

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

negative discount rate

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

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

Batch problem: minimize

$$t_i = d_i - a_i$$

 $a_i$  = arrival time of job i

 $d_i$  = departure time of job i



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

can't start i before  $a_i$ 

Classic metric: mean weighted discounted departure time

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

negative discount rate

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

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

Batch problem: minimize

$$t_i = d_i - a_i$$

 $a_i$  = arrival time of job i

 $d_i$  = departure time of job i



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

can't start i before  $a_i$ 

Classic metric: mean weighted discounted departure time

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

negative discount rate

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

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



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?





Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?







Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?



How do we achieve strong tail optimality?



• Naturally robust to noise



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?





- Naturally robust to noise
- Can adapt to *unknown job sizes* [Harlev et al., 2025]



Why did it take so long to beat FCFS?



Why is achieving strong tail optimality hard?



How does the **Boost** policy family work?

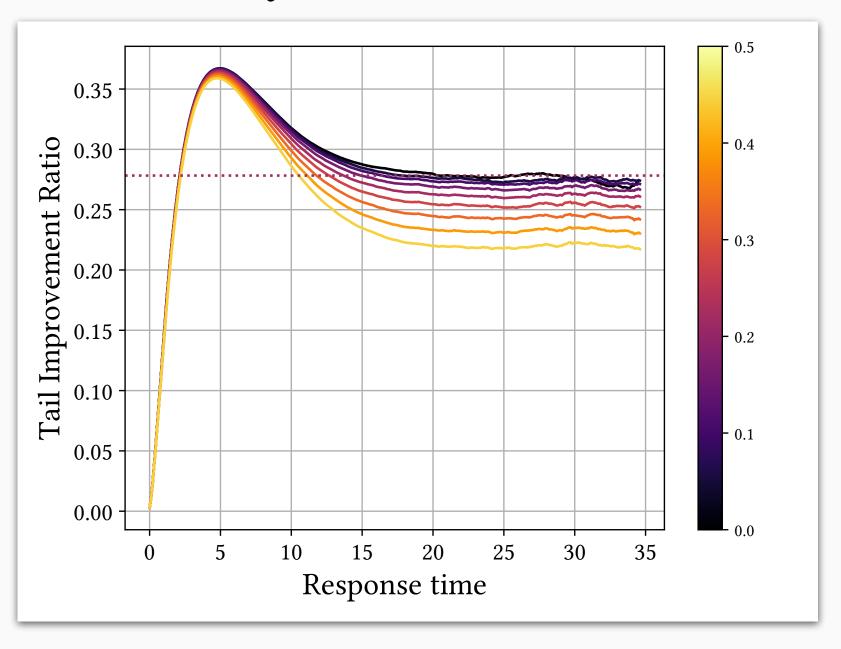




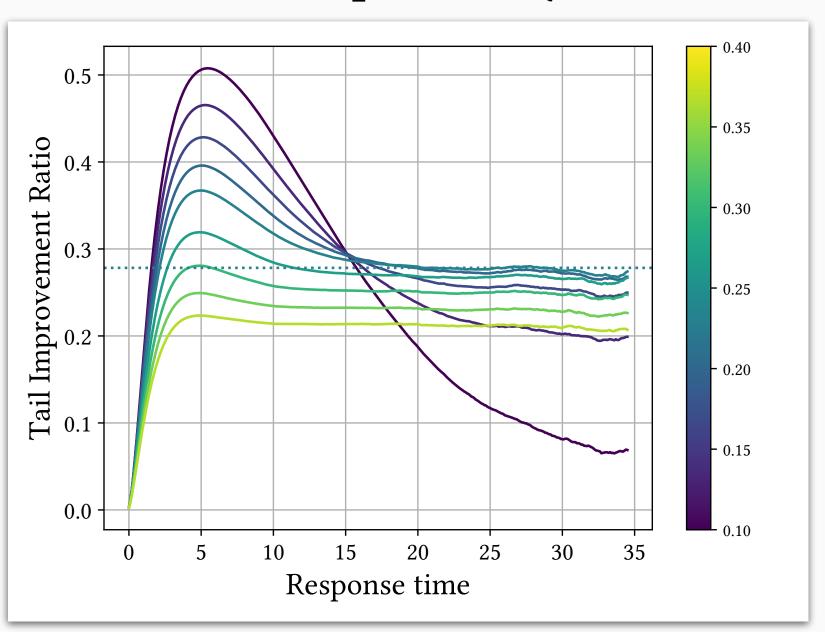
- Naturally robust to noise
- Can adapt to *unknown job sizes* [Harlev et al., 2025]
- Can adapt to multiple servers [Yu et al., 2026]

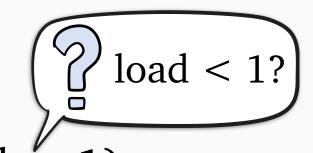
#### Boost is naturally robust

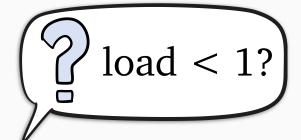
#### Noisy size information

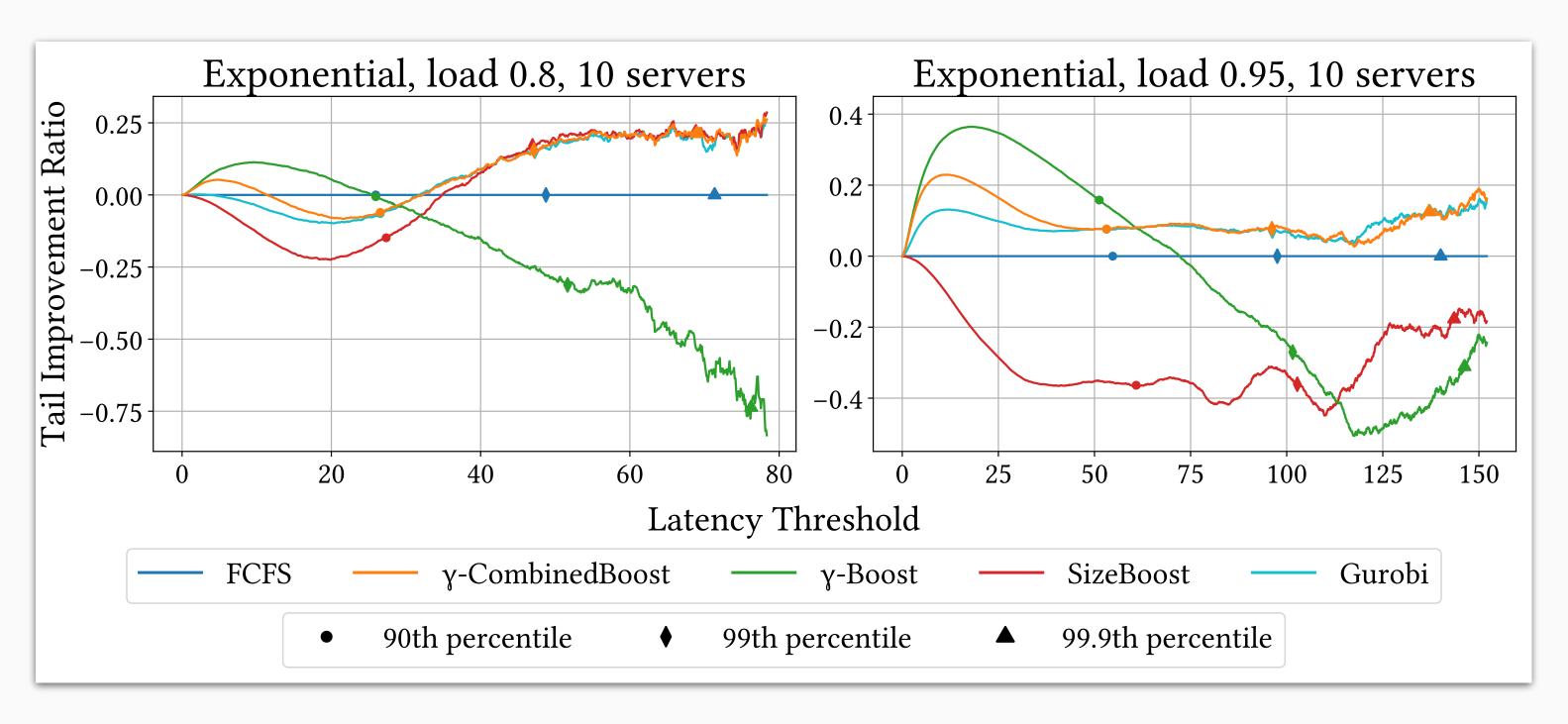


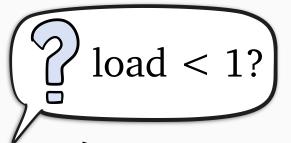
#### Misspecified γ

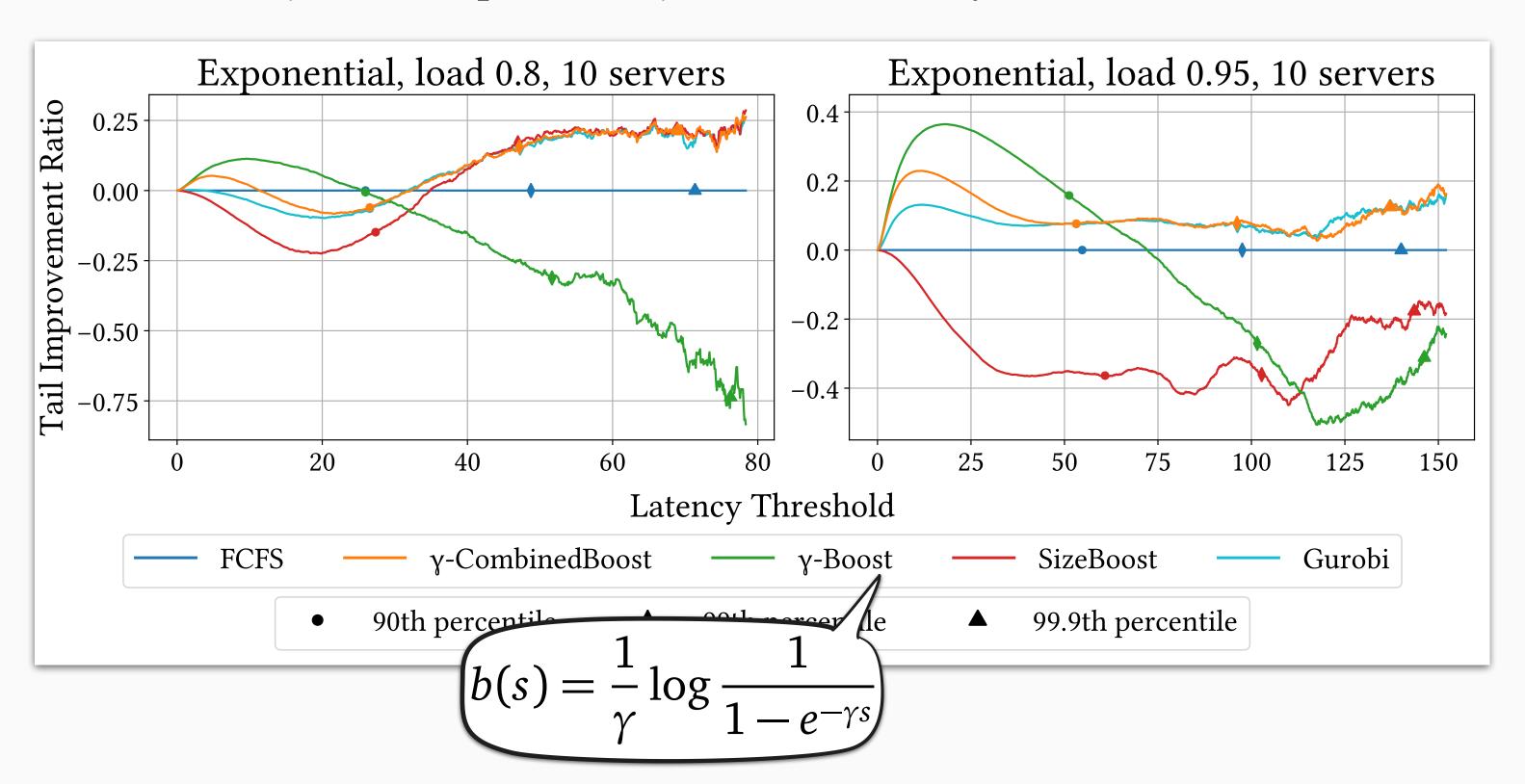


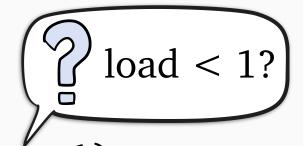


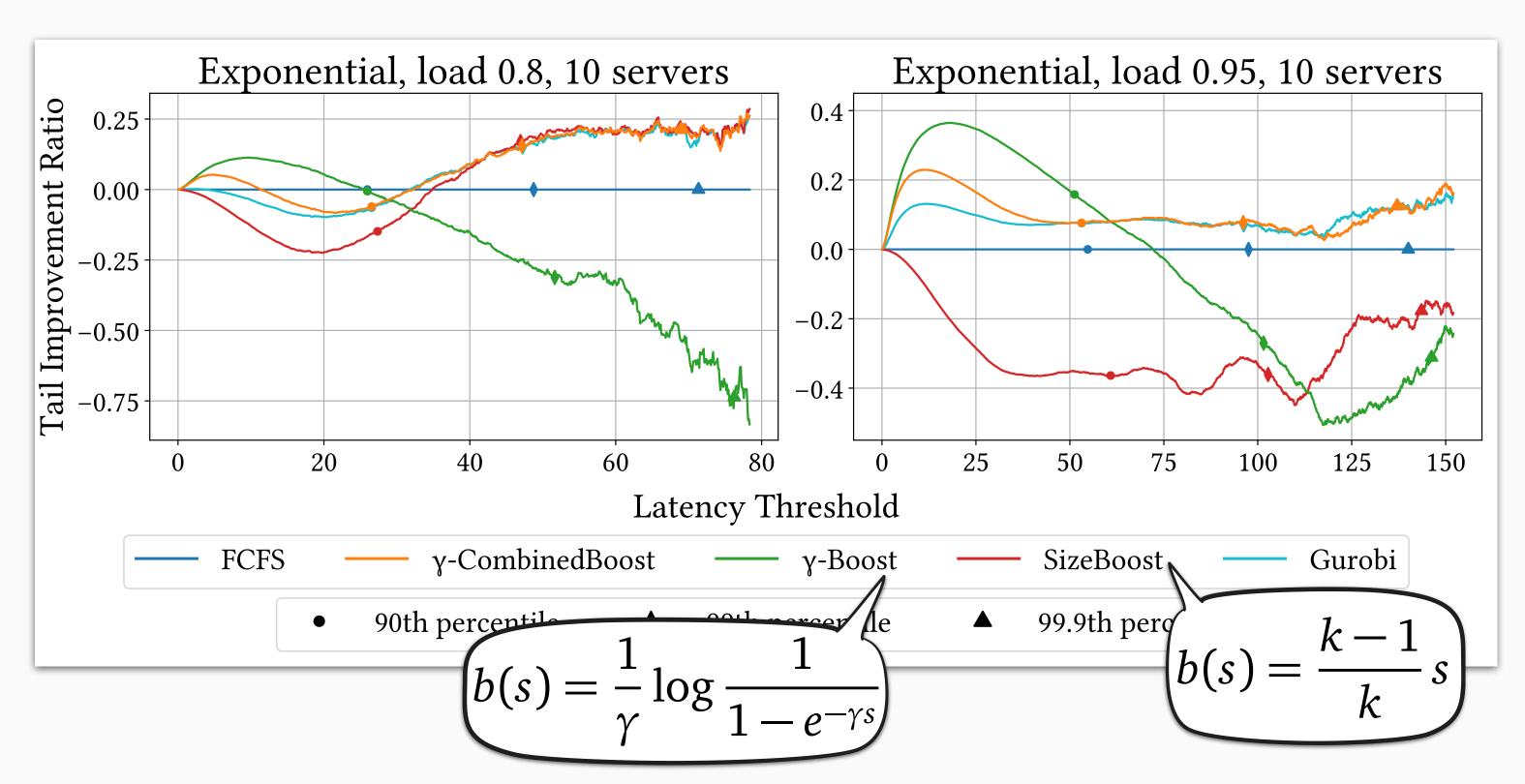


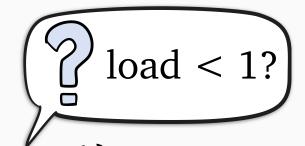


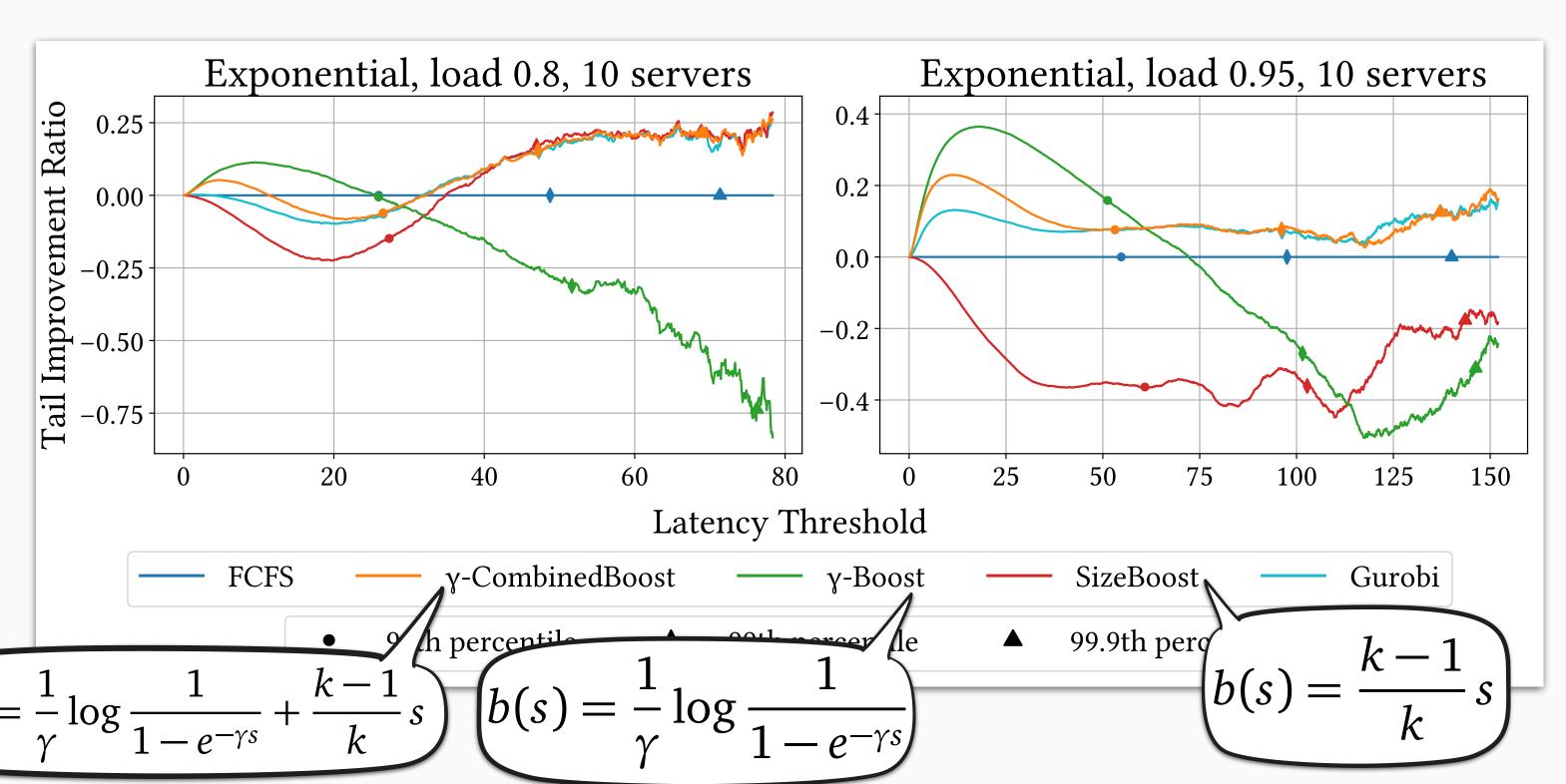














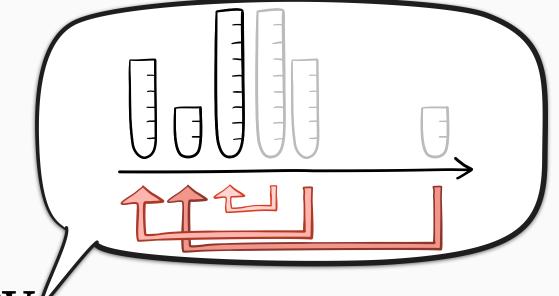
Design the **Boost** scheduling policy



Analyze **Boost**'s performance



Prove Boost is strongly tail-optimal for light-tailed sizes





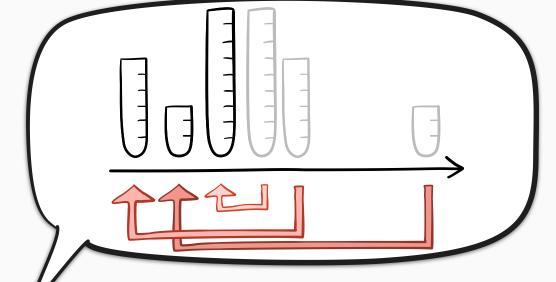
Design the Boost scheduling policy



Analyze Boost's performance



Prove Boost is strongly tail-optimal for light-tailed sizes





Design the **Boost** scheduling policy

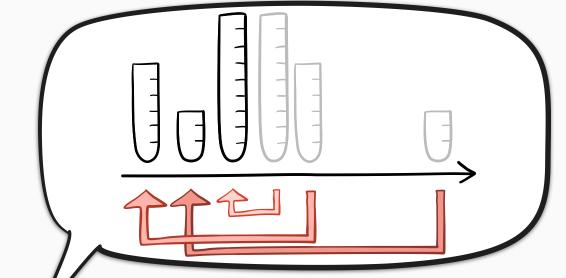


Analyze Boost's performance



Prove Boost is strongly tail-optimal for light-tailed sizes

compute  $C_{\rm Boost}$ 





Design the Boost scheduling policy



Analyze Boost's performance

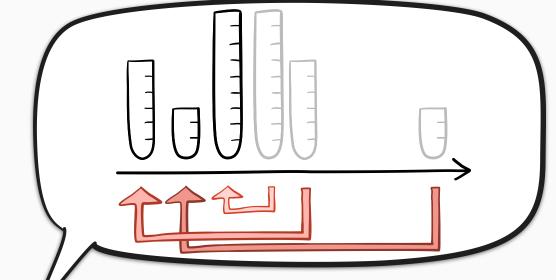


$$y-Boost:$$

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



Prove Boost is strongly tail-optimal for light-tailed sizes





Design the **Boost** scheduling policy



Analyze **Boost**'s performance



$$y-Boost:$$

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

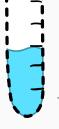


Prove Boost is strongly tail-optimal for light-tailed sizes



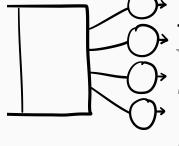
#### Known job sizes

Yu & Scully. Strongly Tail-Optimal Scheduling in the Light-Tailed M/G/1. SIGMETRICS 2024.



#### Unknown job sizes

Harlev, Yu, & Scully. A Gittins Policy for Optimizing Tail Latency. SIGMETRICS 2025.



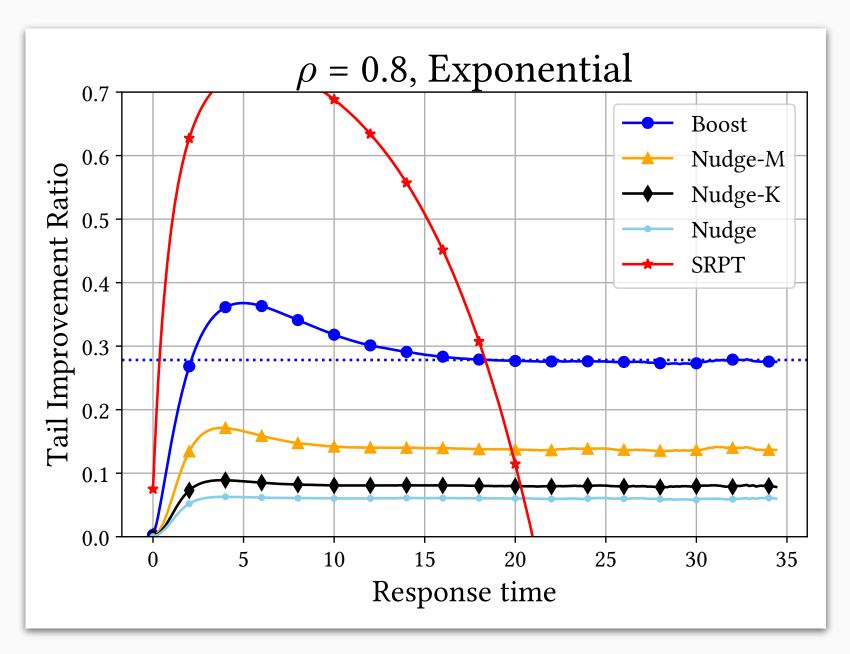
#### Multiple servers

Yu, Harlev, Adakroy, & Scully. *A Tale of Two Traffics: Optimizing Tail Latency in the Light-Tailed M/G/k.* SIGMETRICS 2026.

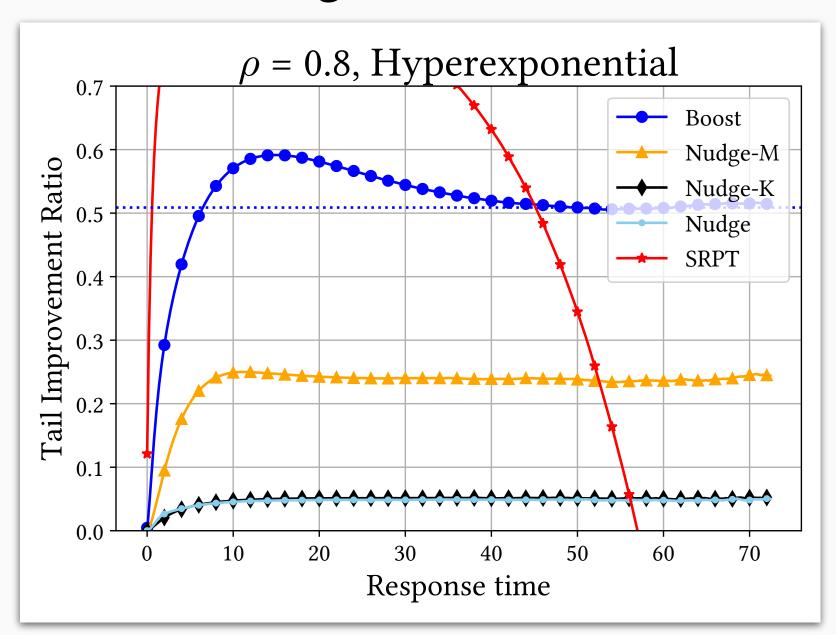
## Bonus slides

#### Impact of job size variance

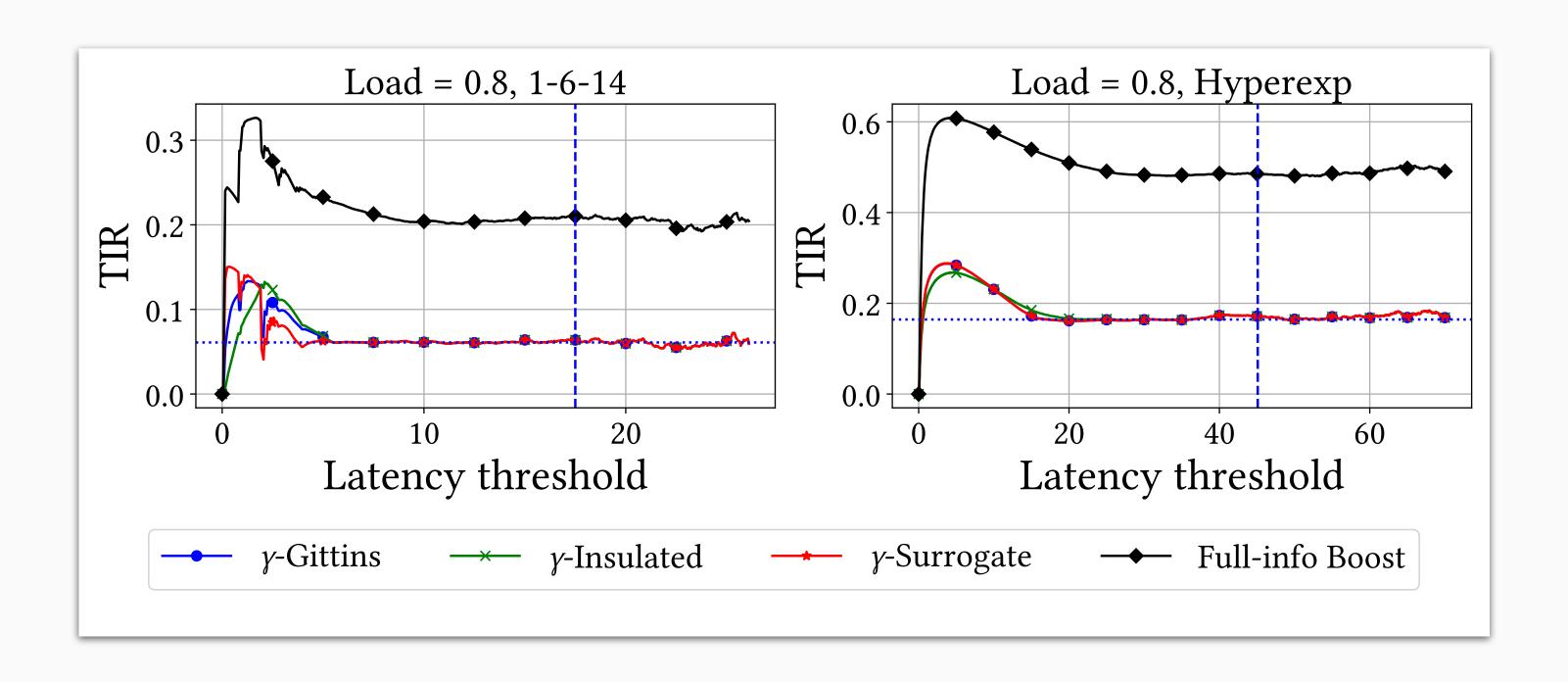
#### Low variance



#### High variance



#### Unknown job sizes



Why is **Boost** weakly tail optimal?

#### Why is **Boost** weakly tail optimal?

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

## Why is **Boost** weakly tail optimal?

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t} \qquad \qquad C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t]$$

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

$$T_{\text{FCFS}} = W + S$$

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

$$T_{\text{FCFS}} = W + S$$
 work

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

$$T_{\text{FCFS}} = W + S$$

$$\text{work}$$

$$C_{\text{FCFS}} = C_W \mathbf{E} [e^{\gamma S}]$$

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

$$T_{\text{FCFS}} = W + S$$

$$Work$$

$$C_{\text{FCFS}} = C_W \mathbf{E}[e^{\gamma S}]$$

$$\lim_{t \to \infty} e^{\gamma t} \mathbf{P}[W > t]$$

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

#### **FCFS**

$$T_{\text{FCFS}} = W + S$$

$$\text{work}$$

$$C_{\text{FCFS}} = C_W \mathbf{E}[e^{\gamma S}]$$

$$\lim_{t \to \infty} e^{\gamma t} \mathbf{P}[W > t]$$

#### **Boost**

$$T_{\text{Boost}} \approx W + S - b(S) + V$$

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

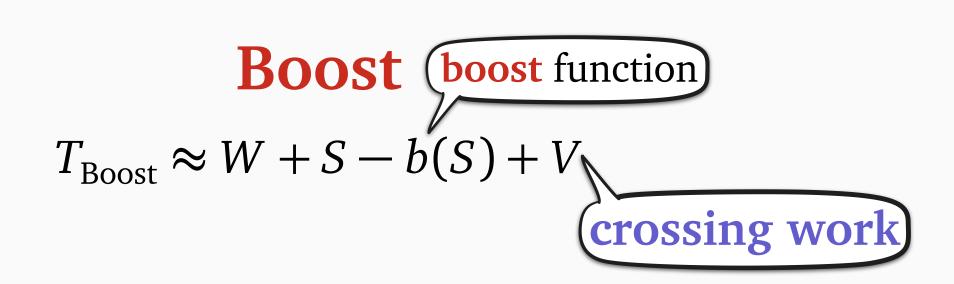
$$T_{\mathrm{FCFS}} = W + S$$
 $\mathrm{work}$ 
 $C_{\mathrm{FCFS}} = C_W \mathbf{E}[e^{\gamma S}]$ 
 $\lim_{t \to \infty} e^{\gamma t} \mathbf{P}[W > t]$ 

Boost boost function
$$T_{\text{Boost}} \approx W + S - b(S) + V$$

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

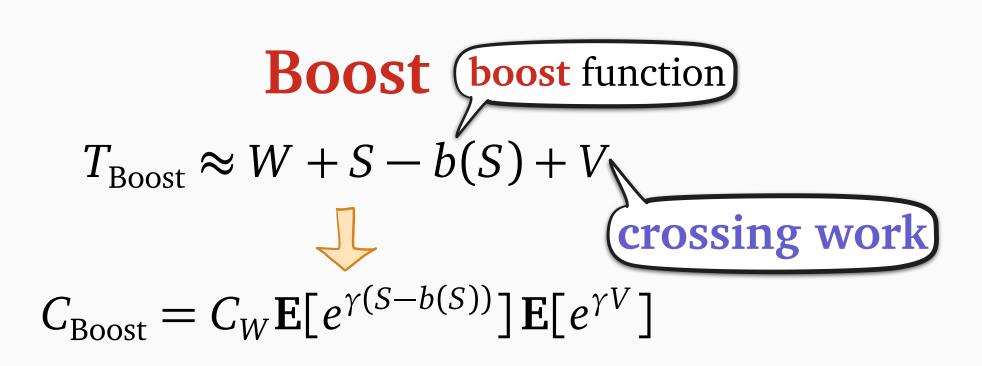
$$T_{ ext{FCFS}} = W + S$$
 $\text{work}$ 
 $C_{ ext{FCFS}} = C_W \mathbf{E}[e^{\gamma S}]$ 
 $\lim_{t \to \infty} e^{\gamma t} \mathbf{P}[W > t]$ 



$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

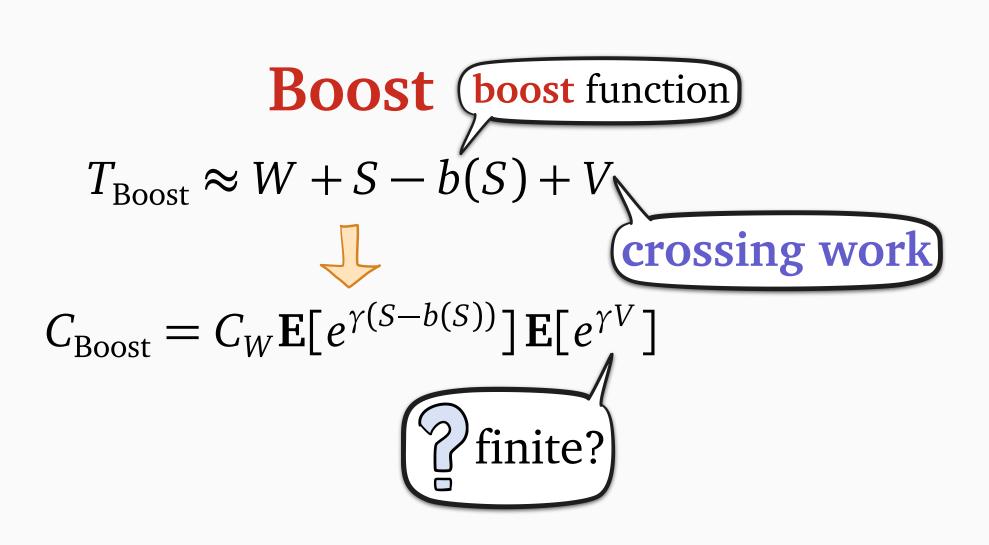
$$T_{ ext{FCFS}} = W + S$$
 $\text{work}$ 
 $C_{ ext{FCFS}} = C_W \mathbf{E}[e^{\gamma S}]$ 
 $\lim_{t \to \infty} e^{\gamma t} \mathbf{P}[W > t]$ 



$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

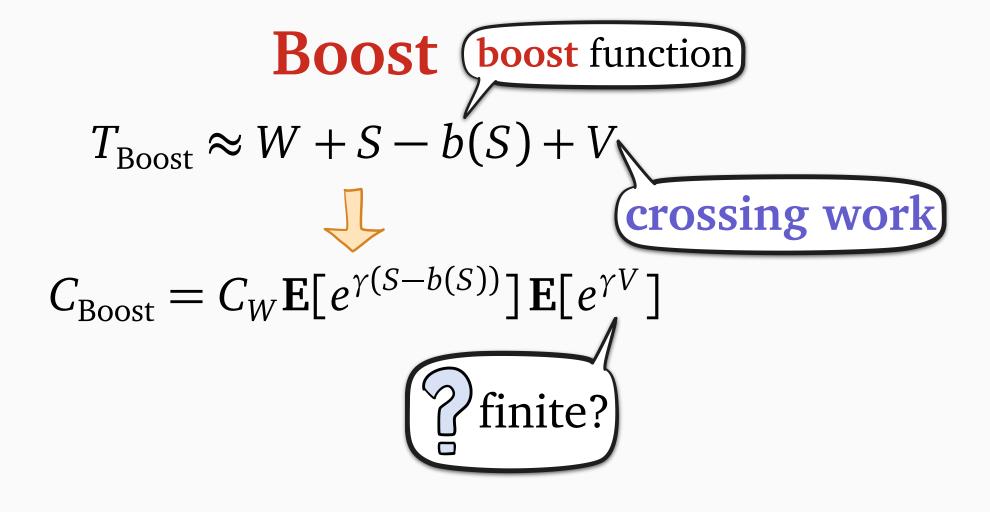
$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

$$T_{ ext{FCFS}} = W + S$$
 $\text{work}$ 
 $C_{ ext{FCFS}} = C_W \mathbf{E}[e^{\gamma S}]$ 
 $\lim_{t \to \infty} e^{\gamma t} \mathbf{P}[W > t]$ 



$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

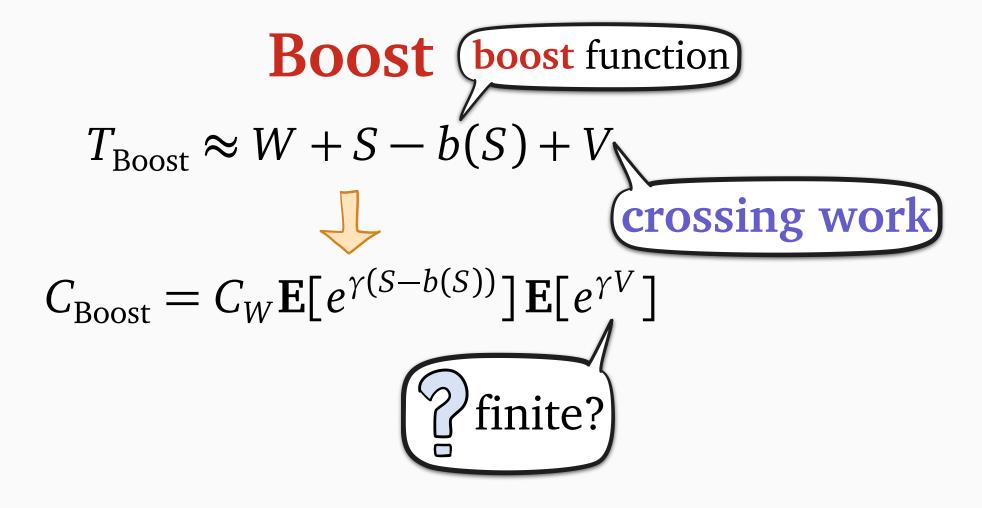
$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem



> time

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem



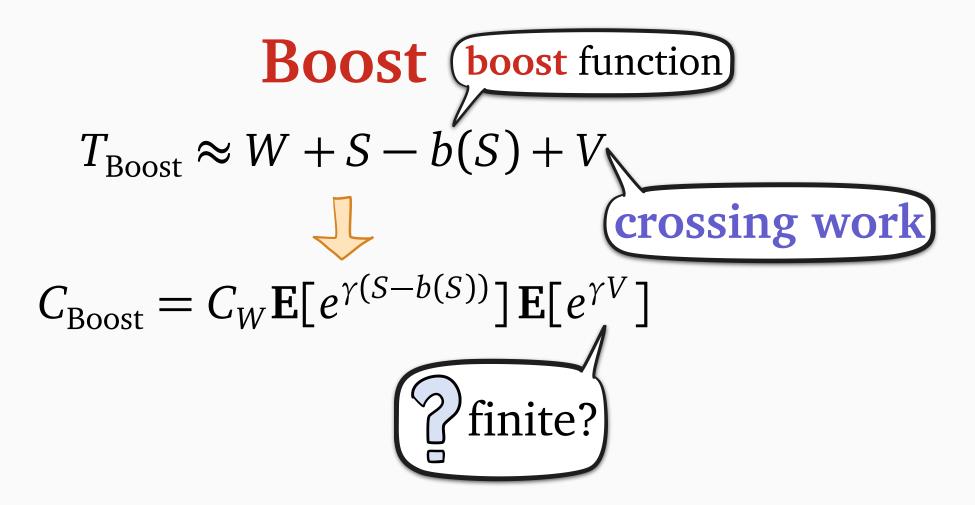
V =crossing work

work that "boosts past" a given time

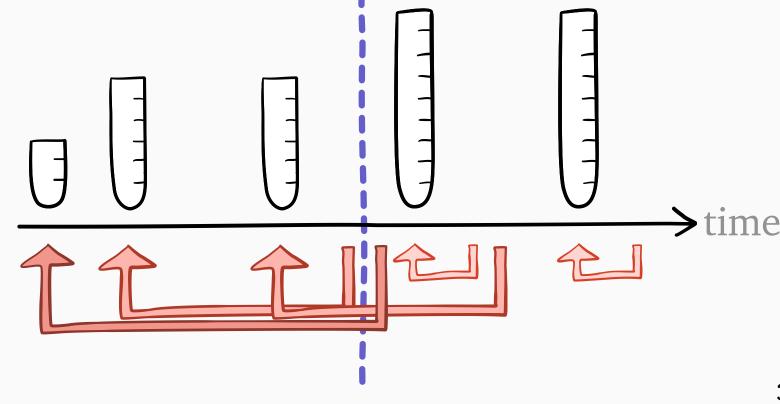
> time

$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

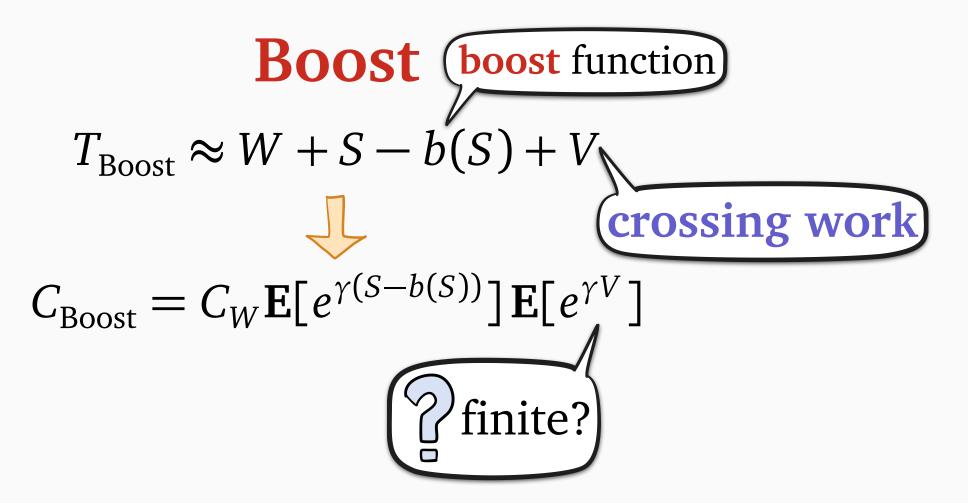


V =crossing work

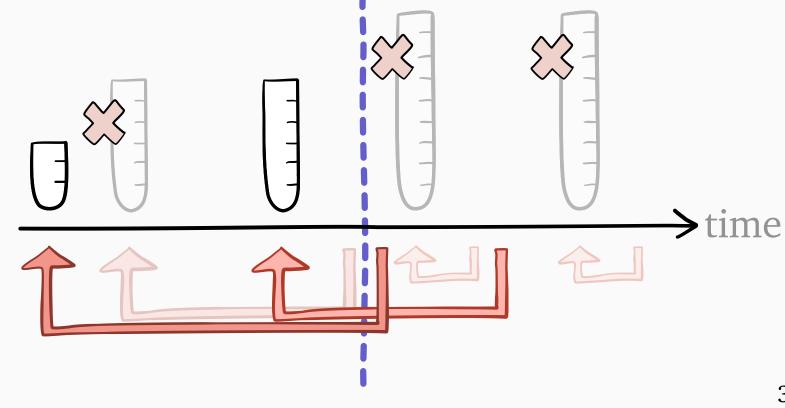


$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

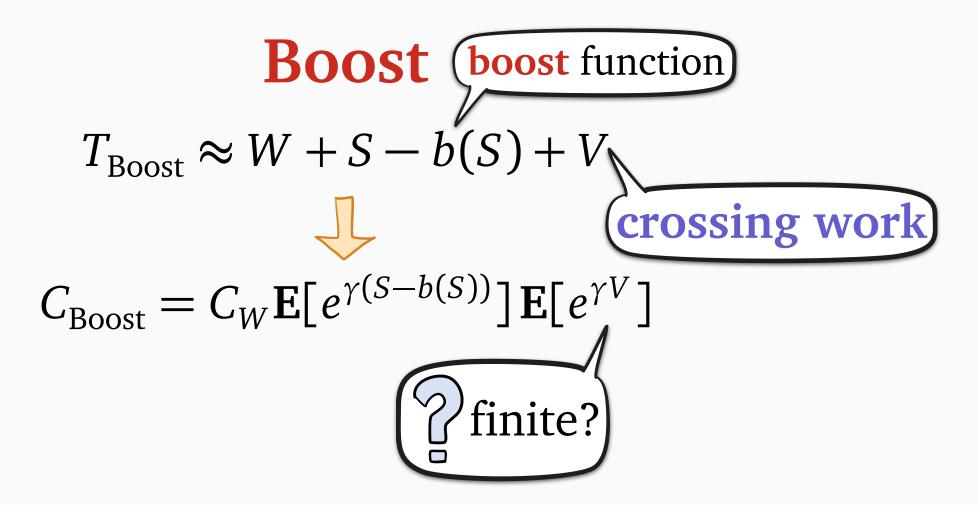


### V =crossing work

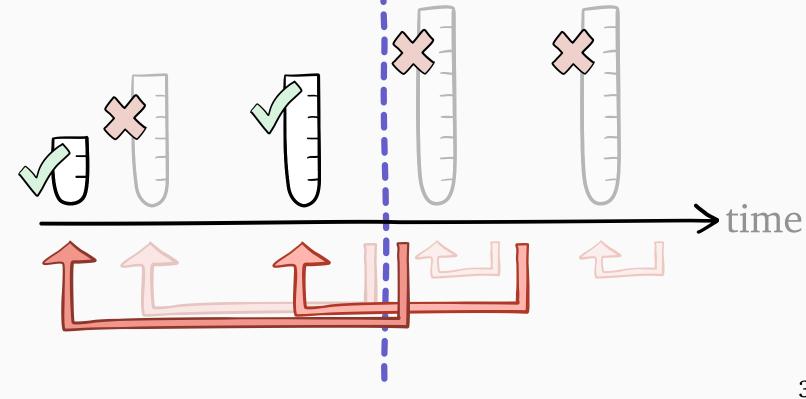


$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem

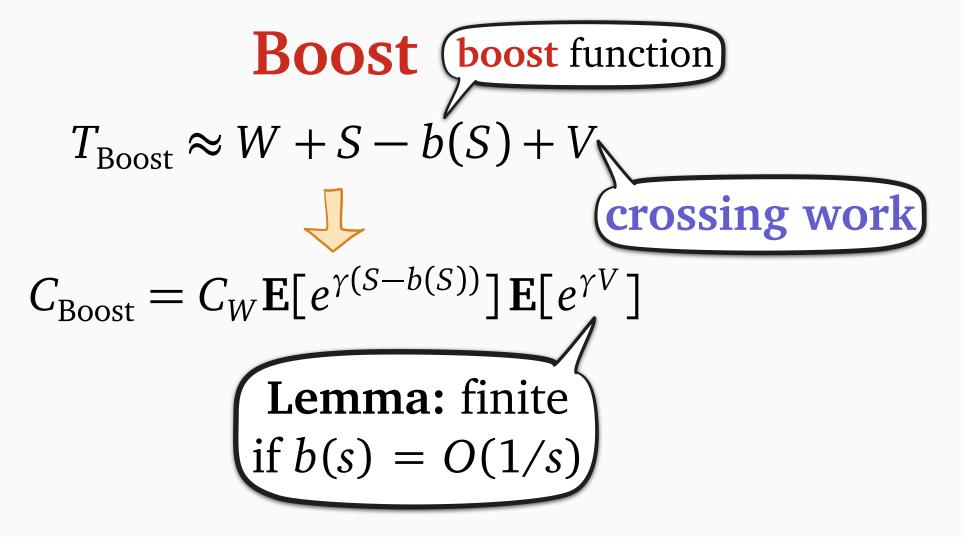


V = crossing work



$$\mathbf{P}[T > t] \sim Ce^{-\gamma t}$$

$$C = \lim_{t \to \infty} e^{\gamma t} \mathbf{P}[T > t] = \lim_{\theta \to \gamma} \frac{\gamma - \theta}{\gamma} \mathbf{E}[e^{\theta T}]$$
final value theorem



$$V = crossing work$$

