# A Gittins Policy for Optimizing Tail Latency

Amit Harlev      Cornell CAM
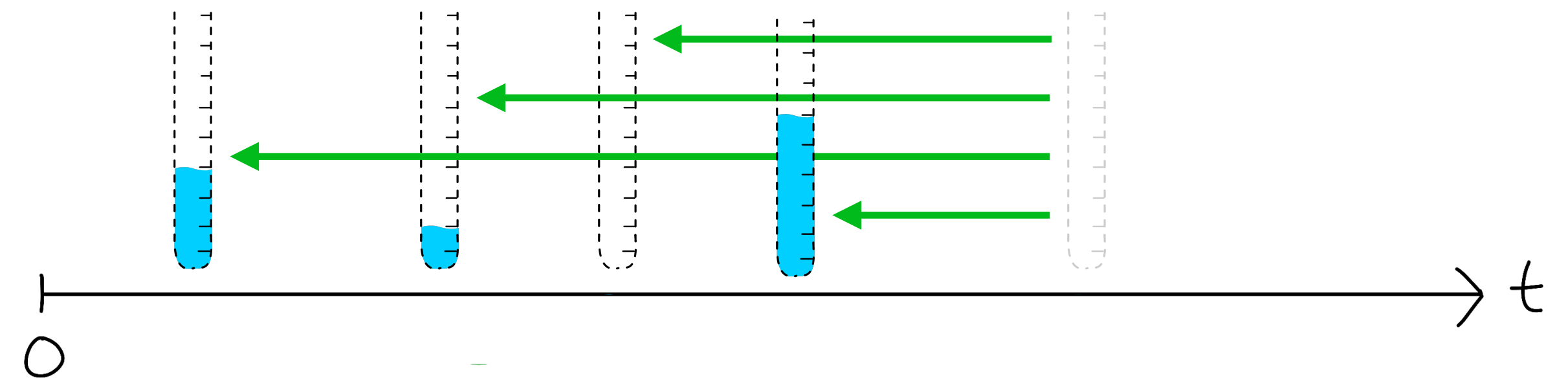
*Joint work with*

George Yu      Cornell ORIE

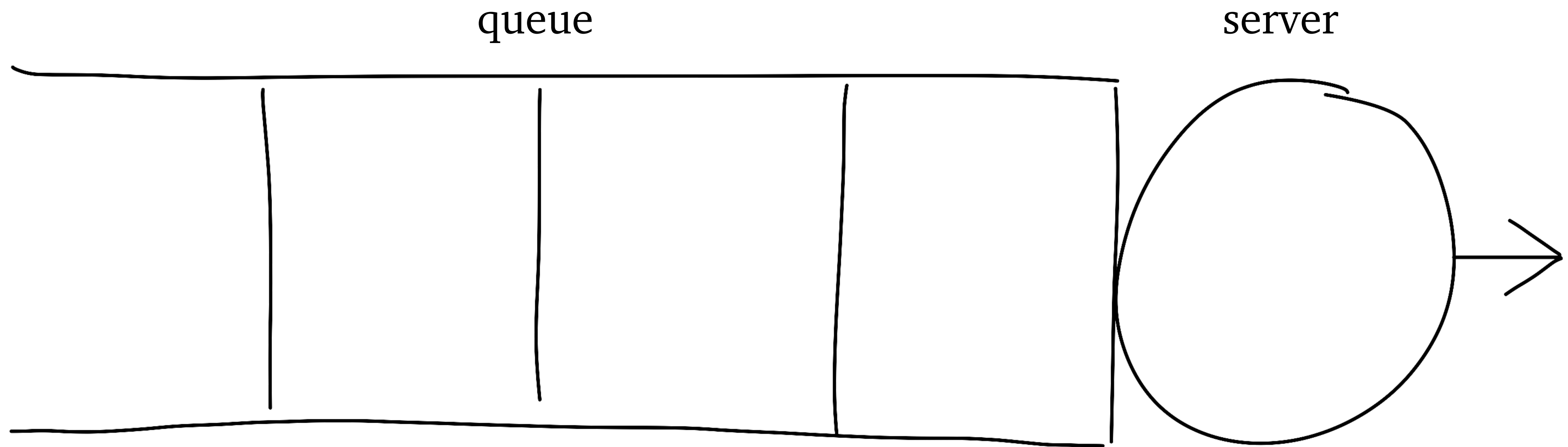Ziv Scully      Cornell ORIE

# How do we minimize delays when job sizes are unknown?

# How do we minimize delays when job sizes are unknown?

(asymptotic) tail latency in single server queue

# Scheduling in the M/G/1

queue

server

# Scheduling in the M/G/1

queue

server

job

size {

# Scheduling in the M/G/1

queue                    server

job

size {

# Scheduling in the M/G/1



queue

server
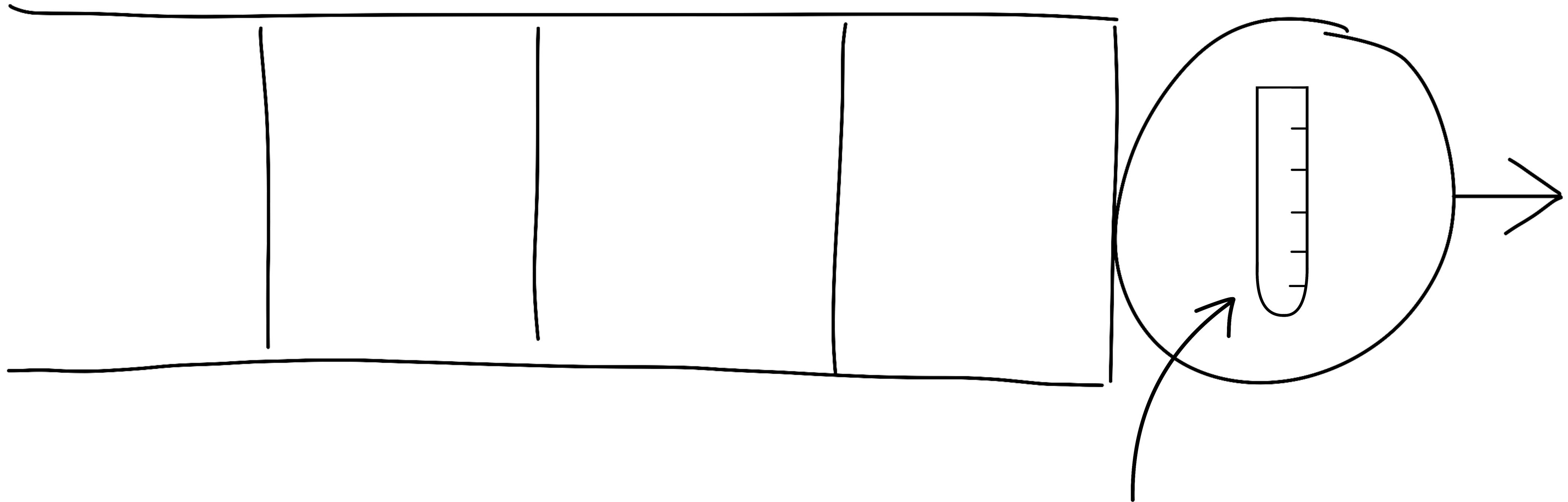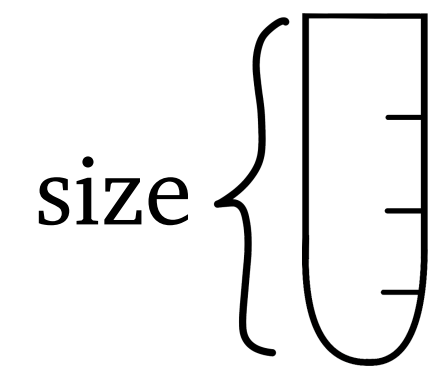
job

size

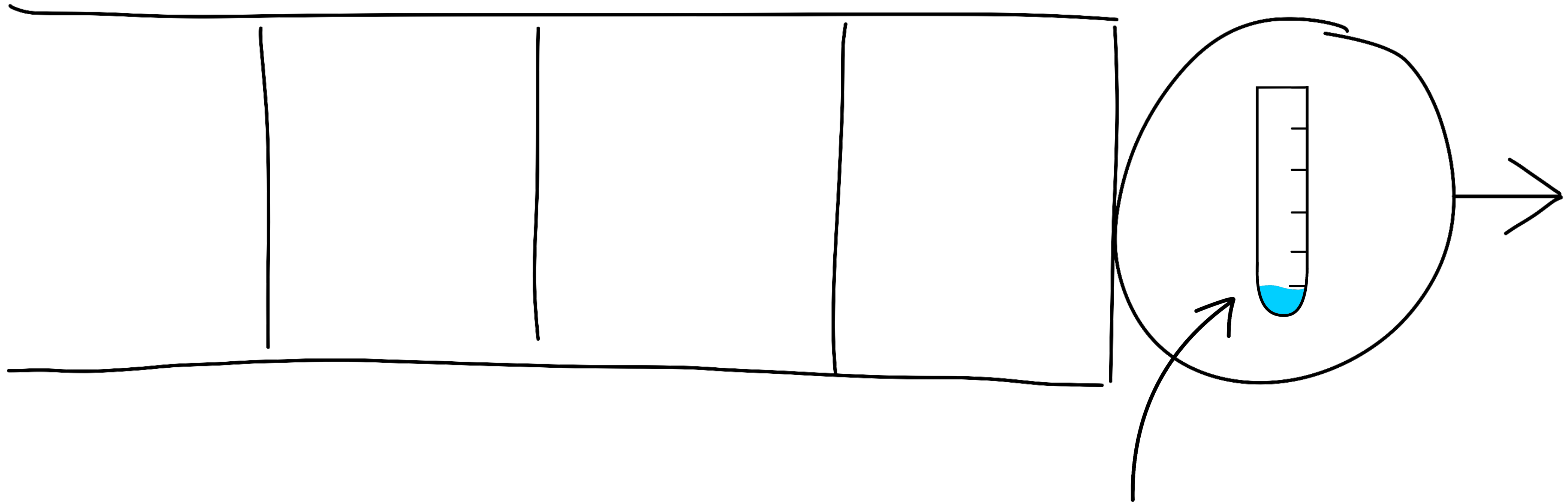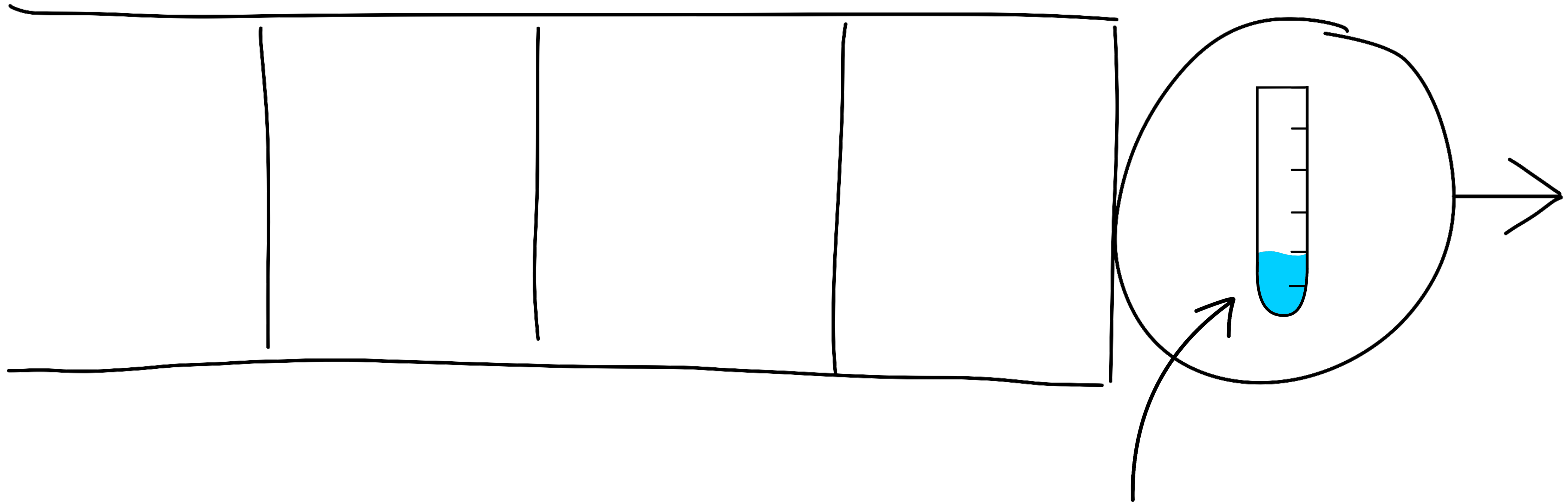# Scheduling in the M/G/1

queue

server

job

size {

# Scheduling in the M/G/1

queue                                                    server



job

size { remaining size

age

# Scheduling in the M/G/1

queue               server

random
arrivals

job

remaining size

size

age

# Scheduling in the M/G/1



queue

server

random arrivals

job

size { remaining size age

# Scheduling in the M/G/1

queue       server

random
arrivals

job

size { remaining size

age

# Scheduling in the M/G/1

queue                                                    server
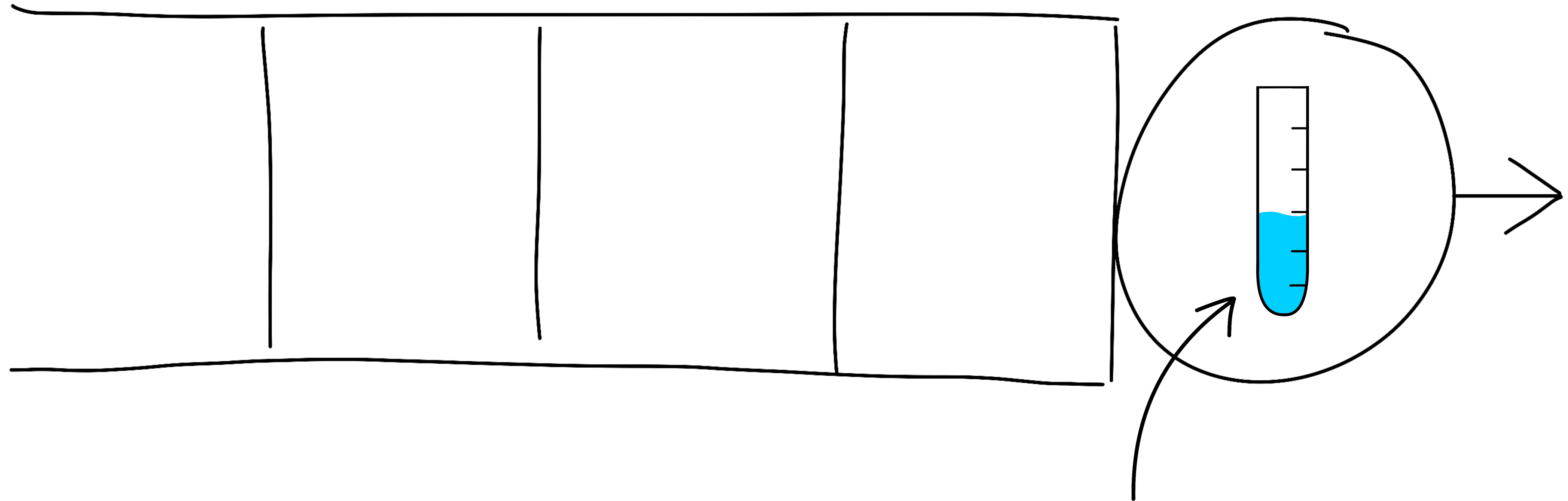
random
arrivals

job



size { remaining size
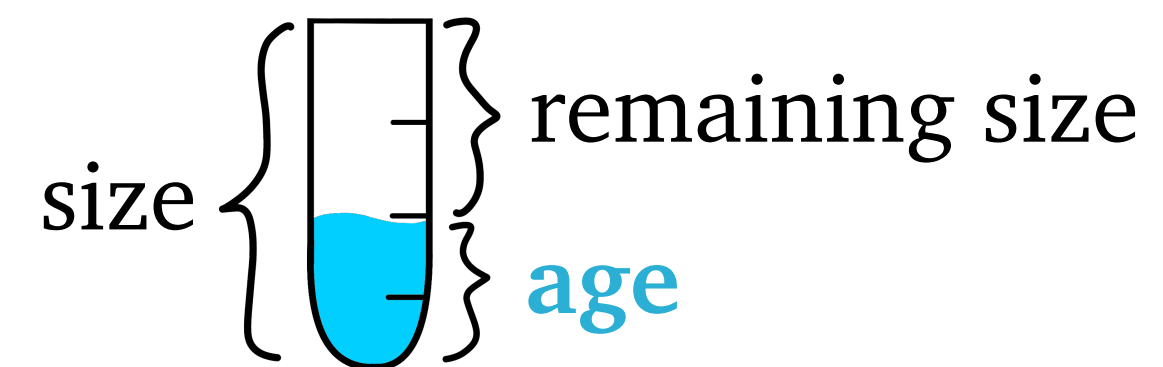     { age

# Scheduling in the M/G/1



queue

server

$S$ = size distribution
$\lambda$ = arrival rate (Poisson)

random arrivals

job

size { remaining size
age

# Scheduling in the M/G/1



Scheduling: In which order should we serve jobs to minimize a desired metric?

queue          server

$T$ = response time

**Scheduling**: In which order should we serve jobs to minimize a desired metric?

queue

server

$T = response\ time$

🎓 **Theory**

queue                                    server

$T$ = response time

**Theory**
---

- mean response time, $\mathbf{E}[T]$

queue                                          server

$T = response\ time$

SRPT          🎓 **Theory**
_____

• mean response time, $\mathbf{E}[T]$

queue                    server

$T$ = response time

**SRPT**          🎓 **Theory**

- mean response time, $\mathbf{E}[T]$
- $\mathbf{E}[T]$ with unknown job sizes

queue server

$T$ = *response time*

SRPT

🎓 **Theory**

- mean response time, $\mathbf{E}[T]$
- $\mathbf{E}[T]$ with unknown job sizes

Gittins

21

queue                                                    server

$T$ = *response time*

SRPT

🎓 **Theory**

- mean response time, $\mathbf{E}[T]$

- $\mathbf{E}[T]$ with unknown job sizes

- weighted $\mathbf{E}[T]$

Gittins

queue                                              server



$T = response\ time$

SRPT          🎓 **Theory**                    ⛑️ **Practice**

- mean response time, $\mathbf{E}[T]$

- $\mathbf{E}[T]$ with unknown job sizes

- weighted $\mathbf{E}[T]$          Gittins

queue                                        server



$T = response\ time$

**SRPT**

🎓 **Theory**

- mean response time, $\mathbf{E}[T]$

- $\mathbf{E}[T]$ with unknown job sizes

- weighted $\mathbf{E}[T]$

**Gittins**

⛑️ **Practice**

- tail latency, $\mathbf{P}[T > t]$ for large $t$

queue                                        server

$T$ = *response time*

**SRPT**

🎓 **Theory**

- mean response time, $\mathbf{E}[T]$
- $\mathbf{E}[T]$ with unknown job sizes
- weighted $\mathbf{E}[T]$

**Gittins**

👷 **Practice**

- tail latency, $\mathbf{P}[T > t]$ for large $t$
- tail latency with unknown job sizes

queue server



$T$ = *response time*

**SRPT**

🎓 **Theory**

- mean response time, $\mathbf{E}[T]$
- $\mathbf{E}[T]$ with unknown job sizes
- weighted $\mathbf{E}[T]$

**Gittins**

⛑️ **Practice**

**hard** to analyze

- tail latency, $\mathbf{P}[T > t]$ for large $t$
- tail latency with unknown job sizes

queue                                                    server



*...sponse time*

**SRPT**          🎓 Th...

This talk: *asymptotic tail latency*
$\mathbf{P}[T > t]$ as $t \to \infty$
with unknown job sizes!

🦺 **Practice**          **hard** to analyze

- mean respo...

- $\mathbf{E}[T]$ with unknown job sizes

- weighted $\mathbf{E}[T]$

- tail latency, $\mathbf{P}[T > t]$ for large $t$

- tail latency with unknown job sizes

**Gittins**

# What does it mean to minimize asymptotic tail latency?

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

**Which policies are weakly optimal?**

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

**Which policies are weakly optimal?**

## Heavy-Tailed Size Distribution

- PS (Processor Sharing)

- LAS (Least Attained Service)

- SRPT (Shortest Remaining Processing Time)

- PLCFS (Preemptive Last Come First Serve)

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

"$P[S > x] \sim \Omega(x^{-\beta})$"

**Which policies are weakly optimal?**

Heavy-Tailed Size Distribution

- PS (Processor Sharing)

- LAS (Least Attained Service)

- SRPT (Shortest Remaining Processing Time)
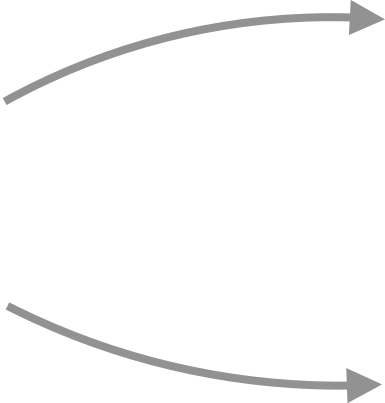
- PLCFS (Preemptive Last Come First Serve)

34

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

**Which policies are weakly optimal?**

*"$P[S > x] \sim \Omega(x^{-\beta})$"*

*"$P[S > x] \sim \mathrm{O}(e^{-\beta x})$"*

## Heavy-Tailed Size Distribution

- PS (Processor Sharing)

- LAS (Least Attained Service)

- SRPT (Shortest Remaining Processing Time)

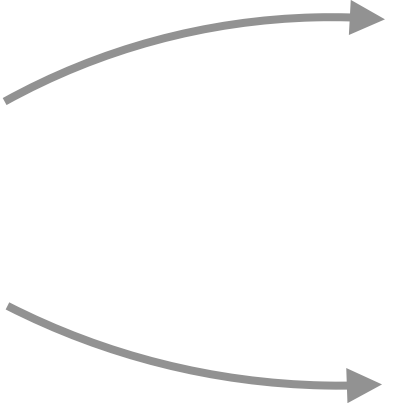- PLCFS (Preemptive Last Come First Serve)

## Light-Tailed Size Distribution

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

**Which policies are weakly optimal?**

"$P[S > x] \sim \Omega(x^{-\beta})$"

"$P[S > x] \sim O(e^{-\beta x})$"

### Heavy-Tailed Size Distribution

- PS (Processor Sharing)

- LAS (Least Attained Service)

- SRPT (Shortest Remaining Processing Time)

- PLCFS (Preemptive Last Come First Serve)

### Light-Tailed Size Distribution

- FCFS (First Come First Serve)

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

**Which policies are weakly optimal?**

*"$P[S > x] \sim \Omega\left(x^{-\beta}\right)$"*

*"$P[S > x] \sim O\left(e^{-\beta x}\right)$"*

### Heavy-Tailed Size Distribution

- PS (Processor Sharing)
- LAS (Least Attained Service)
- SRPT (Shortest Remaining Processing Time)
- PLCFS (Preemptive Last Come First Serve)

### Light-Tailed Size Distribution

- FCFS (First Come First Serve)
- Nudge (Grosof et al. 2021)
- Boost (Yu & Scully 2024)

# What does it mean to minimize asymptotic tail latency?

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

## Which policies are weakly optimal?

### Heavy-Tailed Size Distribution

strongly optimal
- PS (Processor Sharing)
- LAS (Least Attained Service)
- SRPT (Shortest Remaining Processing Time)

- PLCFS (Preemptive Last Come First Serve)

### Light-Tailed Size Distribution

- FCFS (First Come First Serve)

- Nudge (Grosof et al. 2021)

strongly optimal
- Boost (Yu & Scully 2024)

38

# What does it mean to minimize asymptotic tail latency?

## (without job size information)

$$K_\pi = \sup_{\pi^*} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

## Which policies are weakly optimal?

### Heavy-Tailed Size Distribution

strongly optimal

- PS (Processor Sharing)
- LAS (Least Attained Service)
- SRPT (Shortest Remaining Processing Time)
- PLCFS (Preemptive Last Come First Serve)

### Light-Tailed Size Distribution

- FCFS (First Come First Serve)
- Nudge (Grosof et al. 2021)
- Boost (Yu & Scully 2024)

strongly optimal

# What does it mean to minimize asymptotic tail latency?

**(without job size information)**

$$K_\pi = \sup_{\pi^* \in \text{UnknownSize}} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

## Which policies are weakly optimal?

### Heavy-Tailed Size Distribution

strongly optimal
- PS (Processor Sharing)
- LAS (Least Attained Service)
- SRPT (Shortest Remaining Processing Time)
- PLCFS (Preemptive Last Come First Serve)

### Light-Tailed Size Distribution

- FCFS (First Come First Serve)

strongly optimal
- Nudge (Grosof et al. 2021)
- Boost (Yu & Scully 2024)

# What does it mean to minimize asymptotic tail latency?

## (without job size information)

$$K_\pi = \sup_{\pi^* \in \text{UnknownSize}} \lim_{t \to \infty} \frac{\mathbf{P}[T_\pi > t]}{\mathbf{P}[T_{\pi^*} > t]}$$

$K_\pi < \infty$: *weakly optimal*

$K_\pi = 1$: *strongly optimal*

## Which policies are weakly optimal?

### Heavy-Tailed Size Distribution

strongly optimal
- PS (Processor Sharing)
- LAS (Least Attained Service)
- SRPT (Shortest Remaining Processing Time)

- PLCFS (Preemptive Last Come First Serve)

### Light-Tailed Size Distribution

- FCFS (First Come First Serve)

strongly optimal
- Nudge (Grosof et al. 2021)
- Boost (Yu & Scully 2024)

# **Our contribution:** new policy + proof of strong optimality

**GittinsBoost**

# Our contribution: **GittinsBoost** + proof of strong optimality

**How do we generalize Boost to the unknown size setting**

| | known sizes | unknown sizes |
|---|---|---|
| $\lim_{t\to\infty} \mathbf{P}[T > t]$ (light-tailed) | **Boost** $\longrightarrow$ | **GittinsBoost** |

# Our contribution: GittinsBoost + proof of strong optimality



How do we generalize Boost to the unknown size setting

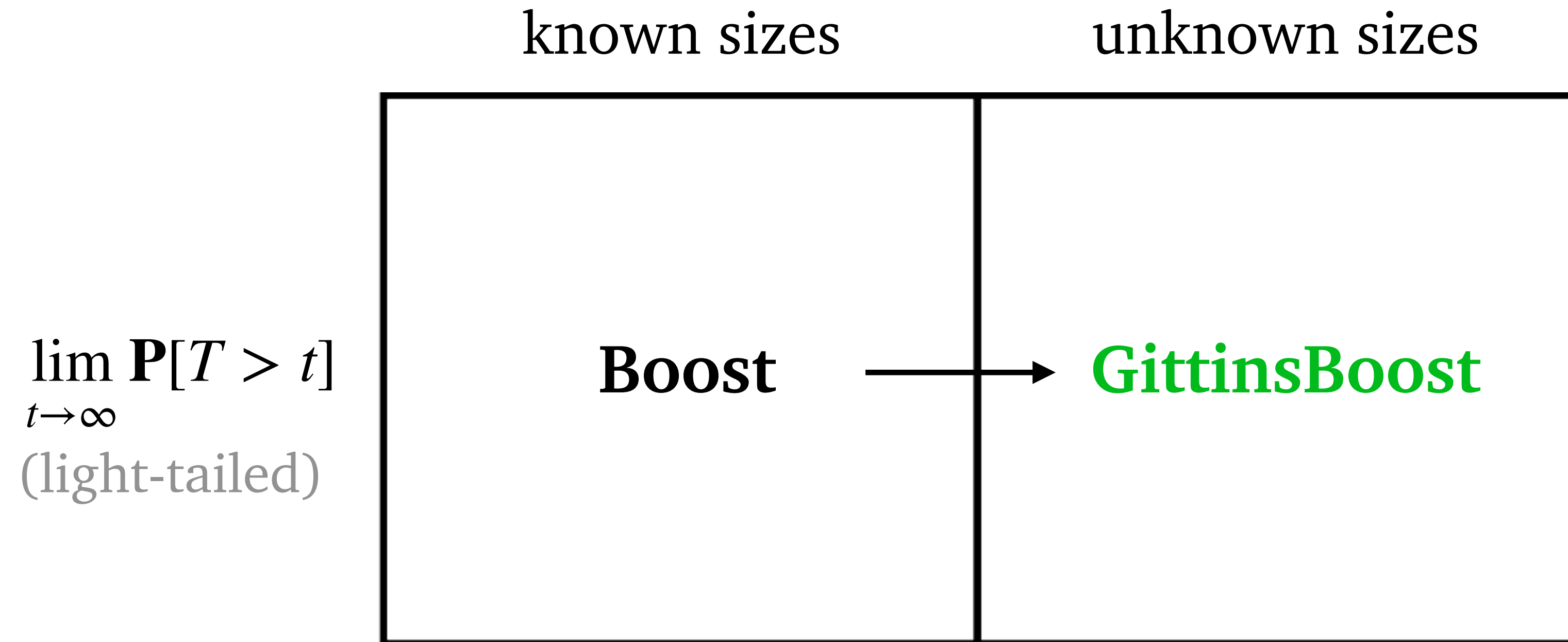| | known sizes | unknown sizes |
|---|---|---|
| $\lim_{t\to\infty} \mathbf{P}[T > t]$ (light-tailed) | **Boost** → | **GittinsBoost** |

What is **Boost**?

# Our contribution: **GittinsBoost** + proof of strong optimality



How do we generalize Boost to the unknown size setting

$$\lim_{t \to \infty} \mathbf{P}[T > t]$$
(light-tailed)

known sizes          unknown sizes

How does scheduling with unknown sizes differ?

**Boost** → **GittinsBoost**

What is **Boost**?

# **Our contribution: GittinsBoost + proof of strong optimality**

# Our contribution: **GittinsBoost** + proof of strong optimality

# Our contribution: **GittinsBoost** + proof of strong optimality

# Our contribution: **GittinsBoost** + proof of strong optimality

# **Our contribution: GittinsBoost + proof of strong optimality**
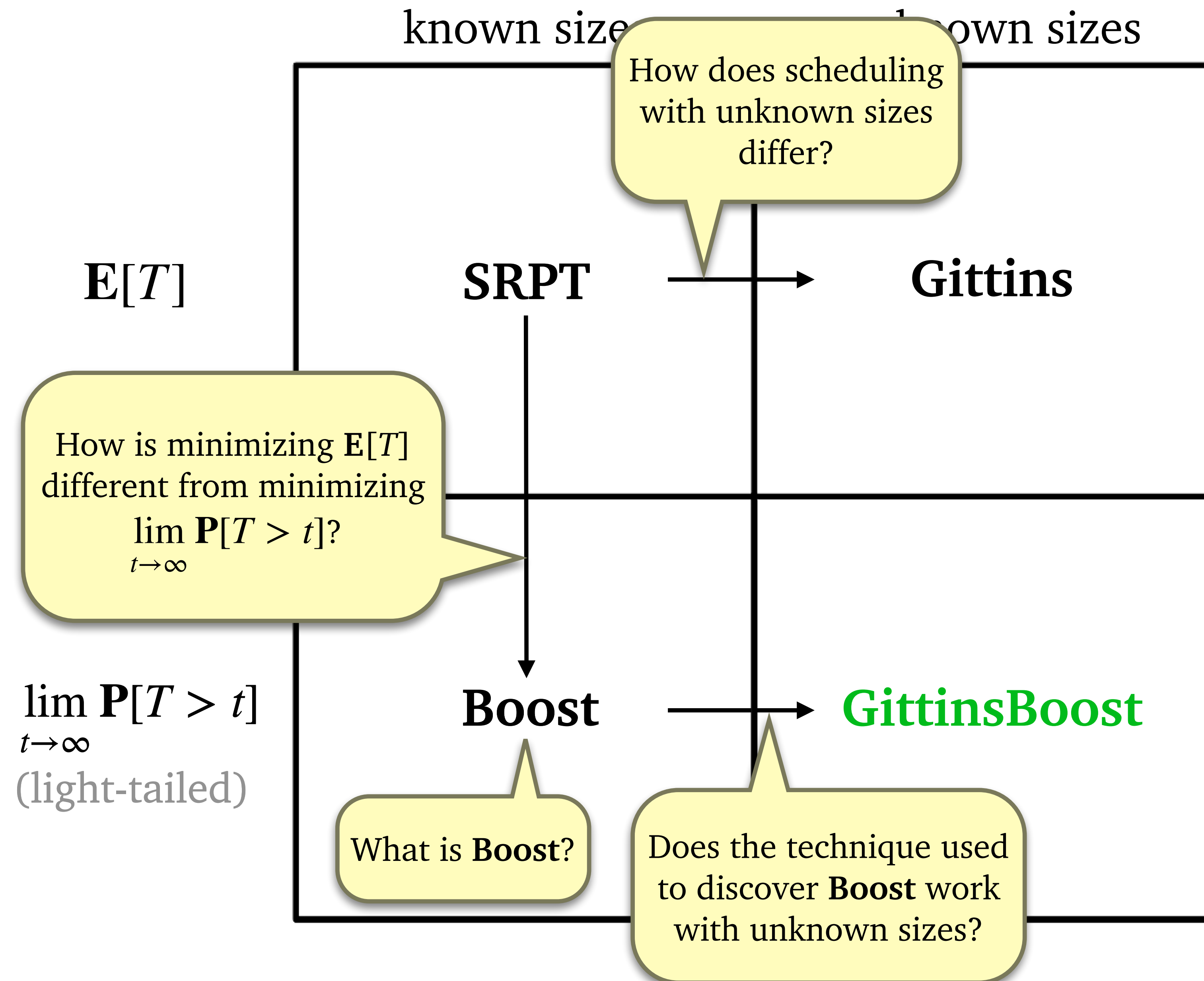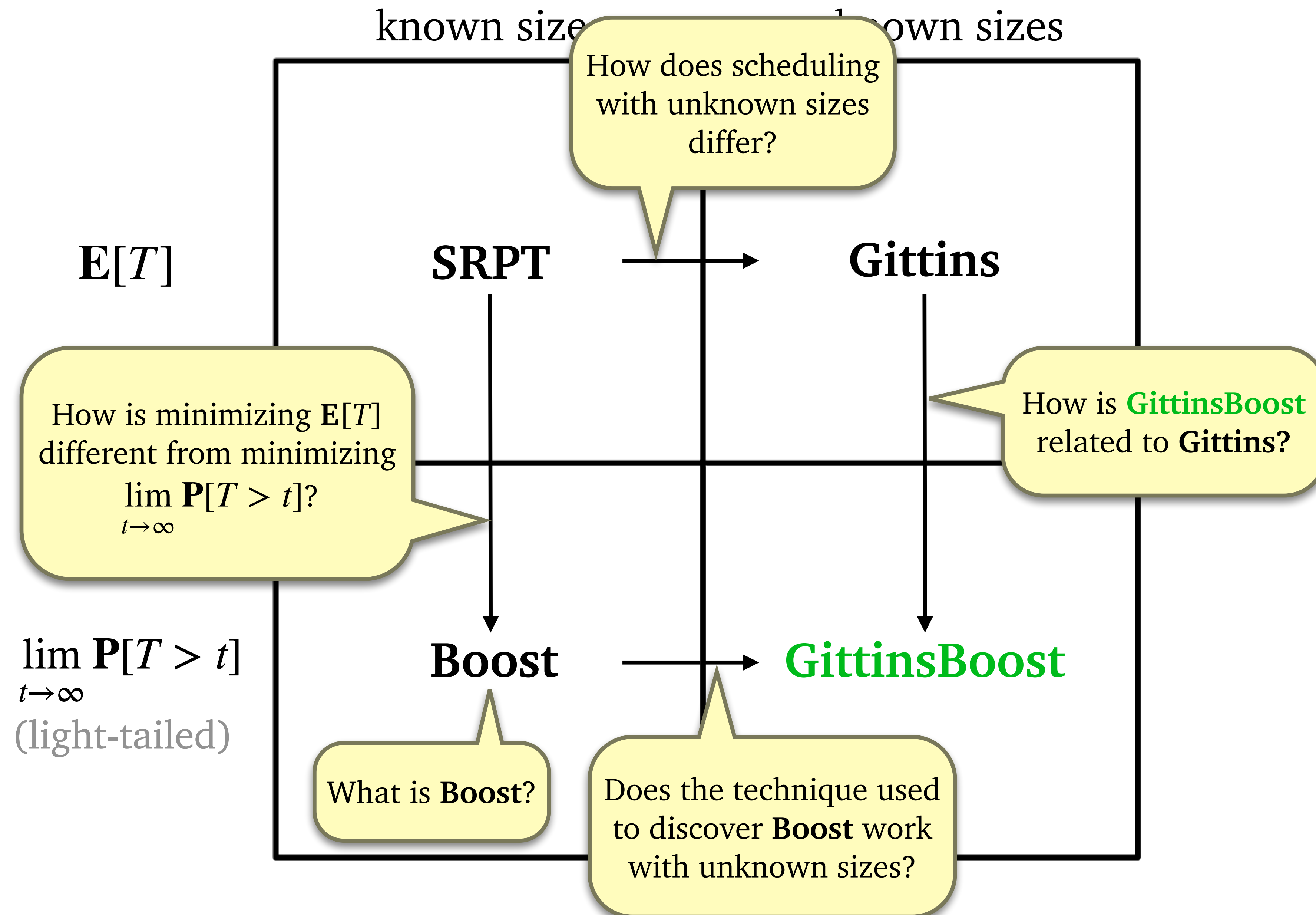
# Optimizing Means vs Optimizing Tails

Minimize $\mathbf{E}[T]$

Minimize $\mathbf{P}[T > t]$

# Optimizing Means vs Optimizing Tails

Minimize $\mathbf{E}[T]$

- Don't want small jobs stuck behind large job

Minimize $\mathbf{P}[T > t]$

# Optimizing Means vs Optimizing Tails

### Minimize $\mathbf{E}[T]$

- Don't want small jobs stuck behind large job

### Minimize $\mathbf{P}[T > t]$

- Don't want small jobs stuck behind large job

# Optimizing Means vs Optimizing Tails

Minimize $\mathbf{E}[T]$

- Don't want small jobs stuck behind large job

Minimize $\mathbf{P}[T>t]$

- Don't want small jobs stuck behind large job

- Don't want large jobs to be overly delayed

# Optimizing Means vs Optimizing Tails

### Minimize $\mathbf{E}[T]$

- Don't want small jobs stuck behind large job

### Minimize $\mathbf{P}[T > t]$

- Don't want small jobs stuck behind large job

- Don't want large jobs to be overly delayed

**Heavy Tails:** basically doesn't matter

# Optimizing Means vs Optimizing Tails

### Minimize $\mathbf{E}[T]$

- Don't want small jobs stuck behind large job

### Minimize $\mathbf{P}[T > t]$

- Don't want small jobs stuck behind large job

- Don't want large jobs to be overly delayed

**Heavy Tails:** basically doesn't matter

**Light Tails:** very important

# Optimizing Means vs Optimizing Tails

Minimize $\mathbf{E}[T]$

- Don't want small jobs stuck behind large job

Minimize $\mathbf{P}[T>t]$

- Don't want small jobs stuck behind large job

- Don't want large jobs to be overly delayed

**Heavy Tails:** basically doesn't matter

**Light Tails:** very important

**Boost:** a way to balance this tradeoff!

# What is Boost?

**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

# What is Boost?

**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function* $b(s) \geq 0$ that maps job size to boost.

# What is Boost?

**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function* $b(s) \geq 0$ that maps job size to boost.

# What is Boost?

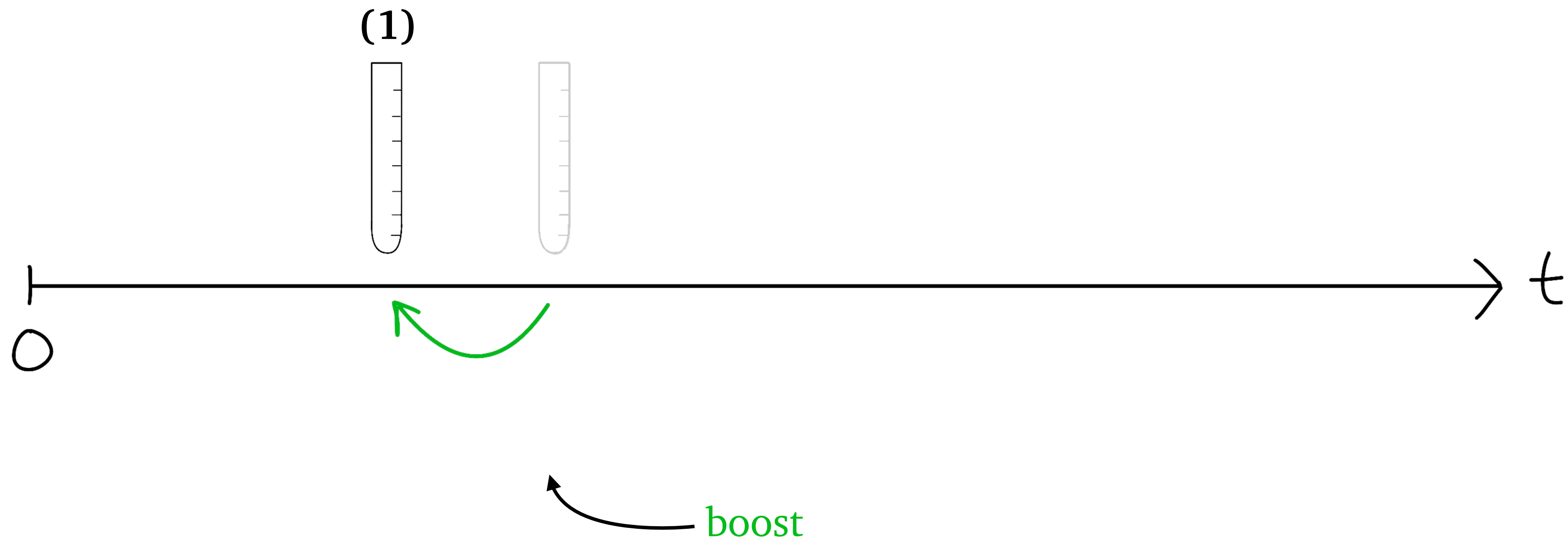**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function $b(s) \geq 0$* that maps job size to boost.



boost

# What is Boost?

**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function* $b(s) \geq 0$ that maps job size to boost.

# What is Boost?

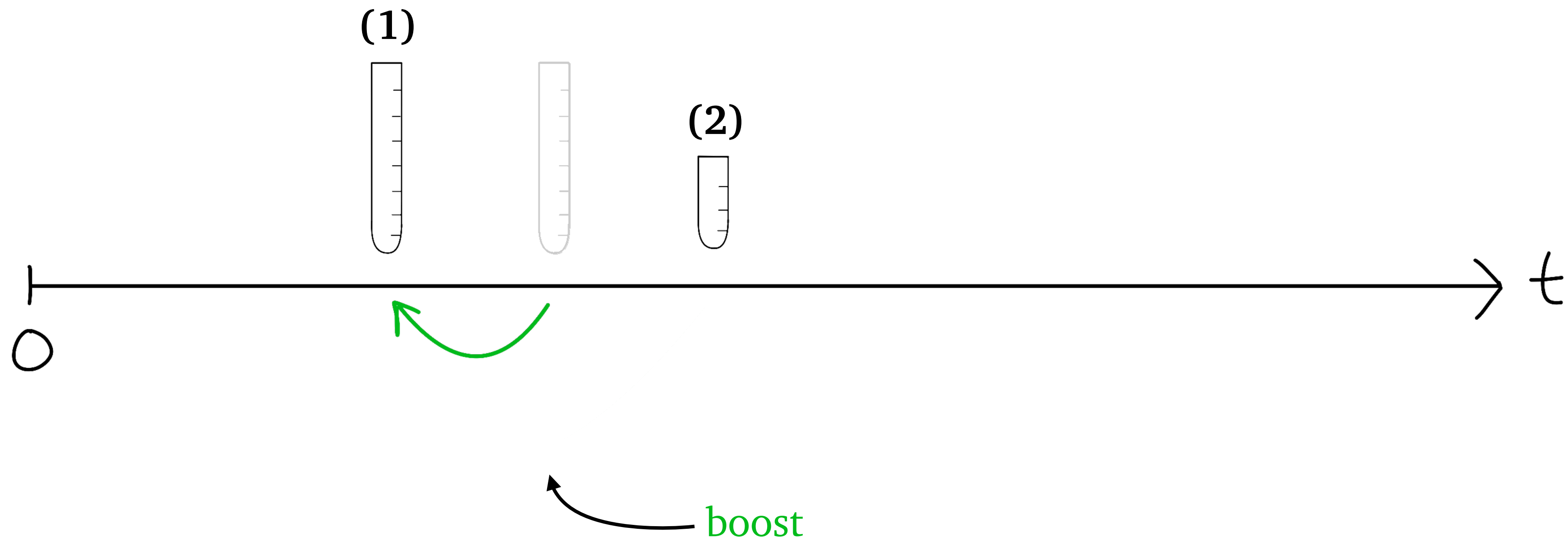**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function $b(s) \geq 0$* that maps job size to boost.

# What is Boost?

**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function* $b(s) \geq 0$ that maps job size to boost.



boost

# What is Boost?
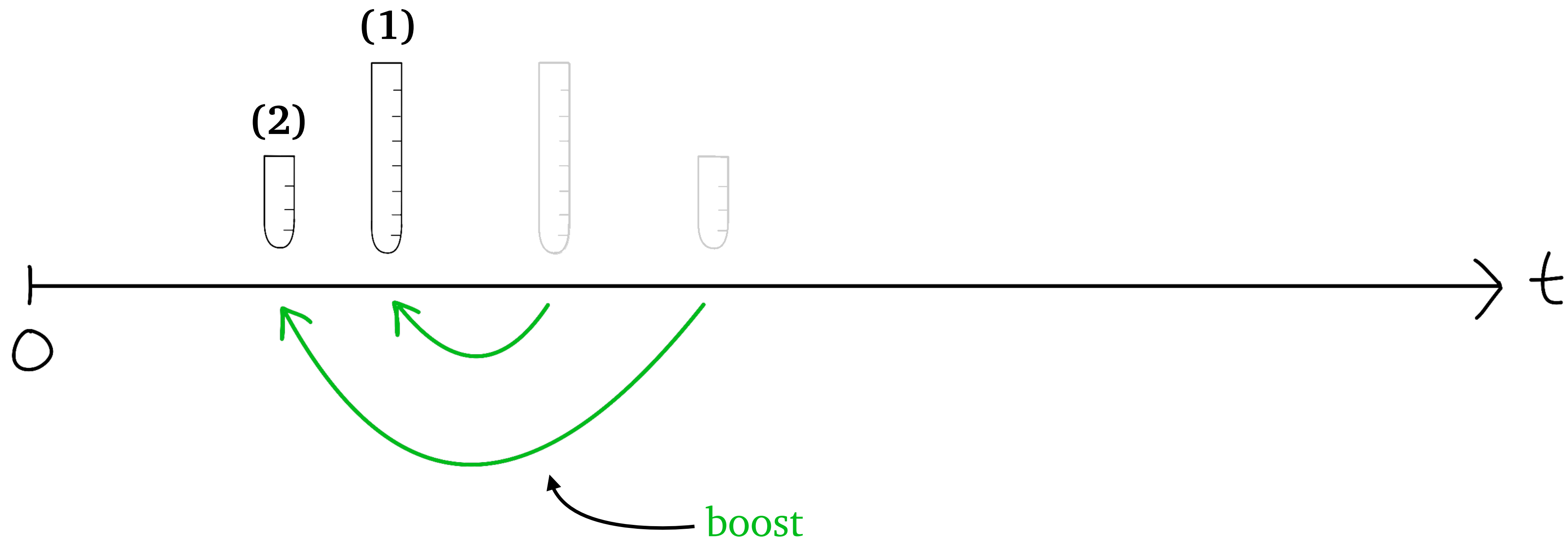
**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function $b(s) \geq 0$* that maps job size to boost.

# What is Boost?

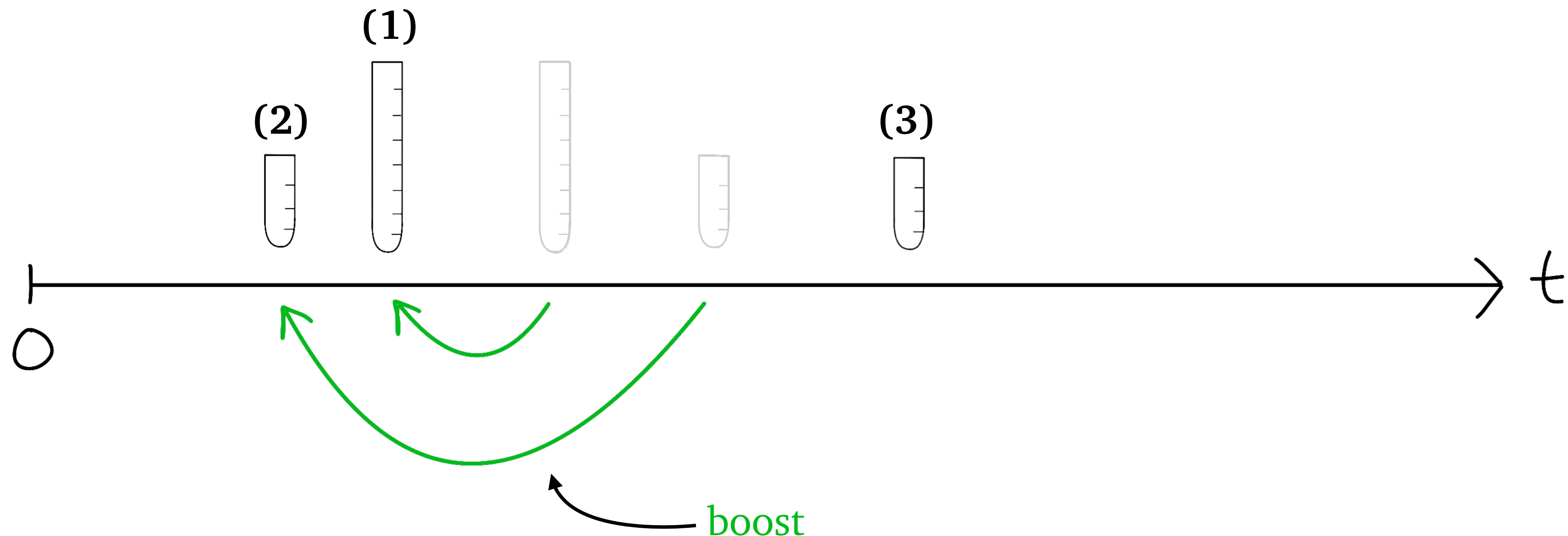**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function $b(s) \geq 0$* that maps job size to boost.

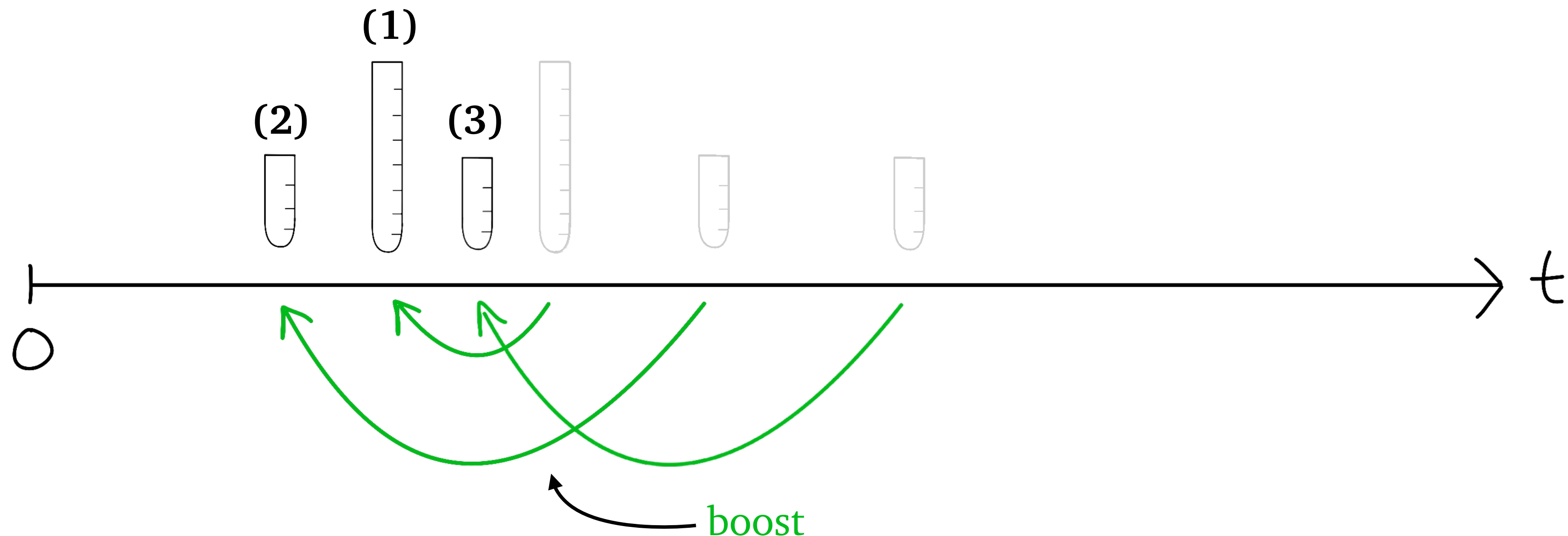**Which boost function minimizes asymptotic tail latency?**

# What is Boost?

**Boost** serves jobs in order of ascending *boosted arrival time:*

boosted arrival time = arrival time - boost

boost is determined by *boost function $b(s) \geq 0$* that maps job size to boost.
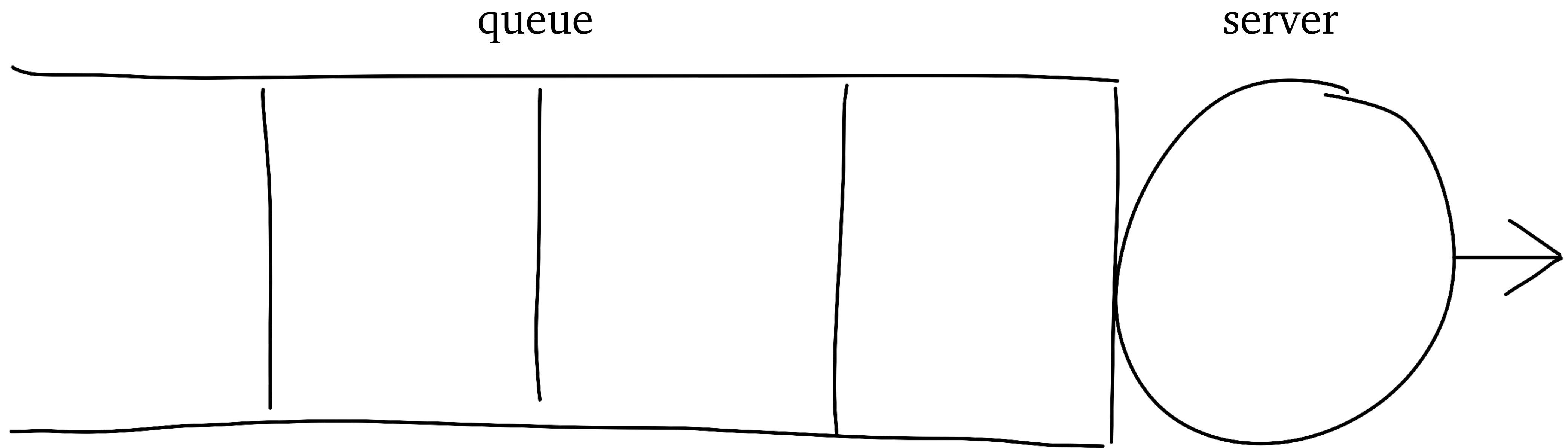
**Which boost function minimizes asymptotic tail latency?**
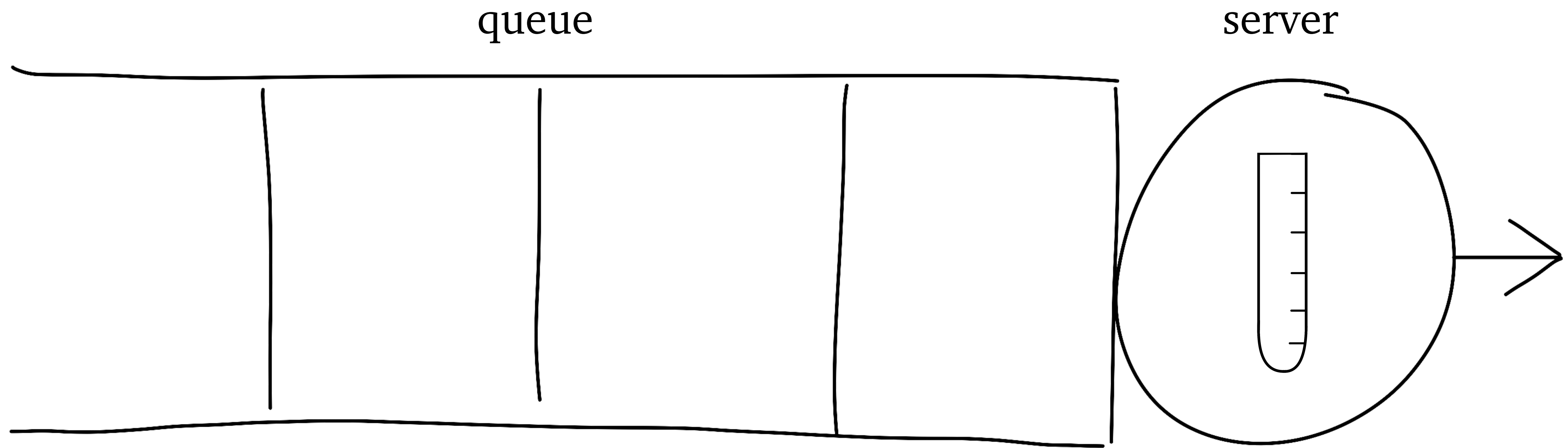
choosing:

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

results in strongly optimal policy (Yu & Scully, 2024)

# Scheduling with unknown job sizes

queue                                    server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue

server

# Scheduling with unknown job sizes

queue                                    server



**state** = (size, age, arrival time)

# Scheduling with unknown job sizes

queue                    server

**state** = (size, age, arrival time)

**SRPT**: order by [size - age]

# Scheduling with unknown job sizes

queue                                                    server

**state** = (size, age, arrival time)

**SRPT**: order by [size - age]

**state** = (age, arrival time)

# Scheduling with unknown job sizes

queue                                          server

**state** = (size, age, arrival time)          **state** = (age, arrival time)
――――――――――――――――――――――               ――――――――――――――――――――
**SRPT**: order by [size - age]                **Gittins**: order by **rank**(age)

rank ← smaller is better

age

# Scheduling with unknown job sizes

queue                                           server

**state** = (size, age, arrival time)        **state** = (age, ~~arrival time~~)

**SRPT**: order by [size - age]             **Gittins**: order by **rank**(age)

proxy for [size - age],
depends on size distribution

smaller is
better

age

# Scheduling with unknown job sizes

queue                                    server



**state** = (size, age, arrival time)

**SRPT**: order by [size - age]

**Boost**: order by [A - *b(s)]*

**state** = (age, arrival time)

**Gittins**: order by **rank**(age)

proxy for [size - age],
depends on size distribution

smaller is
better

age

# Scheduling with unknown job sizes

queue                                        server



proxy for [size - age],
depends on size distribution

**state** = (size, age, arrival time)        **state** = (age, arrival time)

**SRPT**: order by [size - age]              **Gittins**: order by **rank**(age)

**Boost**: order by [A - *b(s)]*             💡 **GittinsBoost**:

smaller is
better

order by [A - *b(*age)]

age

# Scheduling with unknown job sizes

queue

server

**state** = (size, age, arrival time)

**SRPT**: order by [size - age]

**Boost**: order by [A - *b(s)]*

**state** = (age, arrival time)

**Gittins**: order by **rank**(age)

💡 **GittinsBoost**:

order by [A - *b*(age)]

proxy for [size - age],
depends on size distribution

smaller is better

Proxy for *b(s)*

age

# The **GittinsBoost** policy

defined by a *boost function $b(x) \geq 0$* that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function $b(x) \geq 0$* that maps **age** to boost

# The **GittinsBoost** policy

defined by a *boost function* $b(x) \geq 0$ that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function* $b(x) \geq 0$ that maps **age** to boost

# The **GittinsBoost** policy

defined by a *boost function* $b(x) \geq 0$ that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function b(x)* $\geq 0$ that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function $b(x) \geq 0$* that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function* $b(x) \geq 0$ that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function* $b(x) \geq 0$ that maps **age** to boost

# The GittinsBoost policy

defined by a *boost function $b(x) \geq 0$* that maps **age** to boost
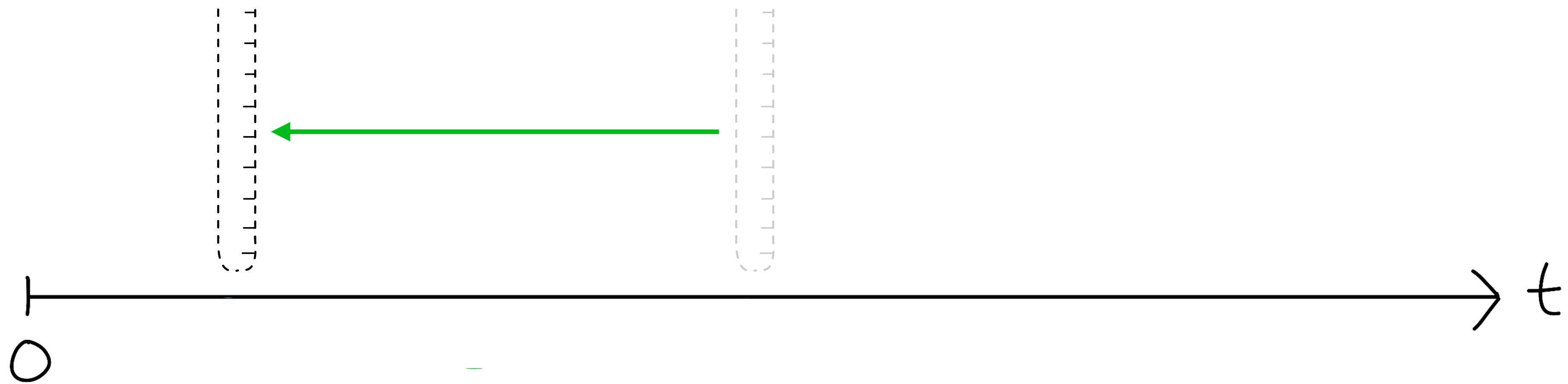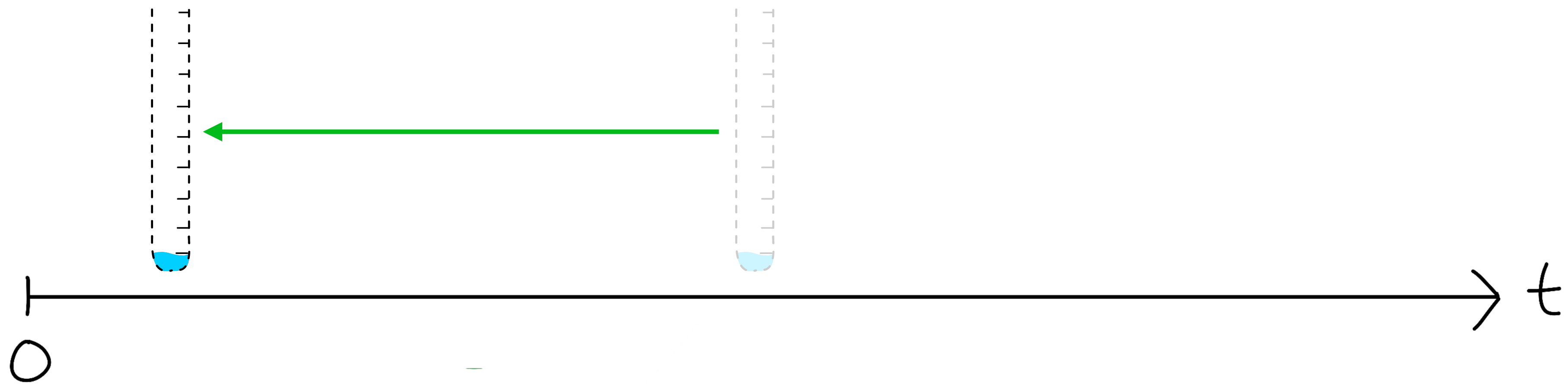


job completes

# The **GittinsBoost** policy

defined by a *boost function* $b(x) \geq 0$ that maps **age** to boost

# The GittinsBoost policy

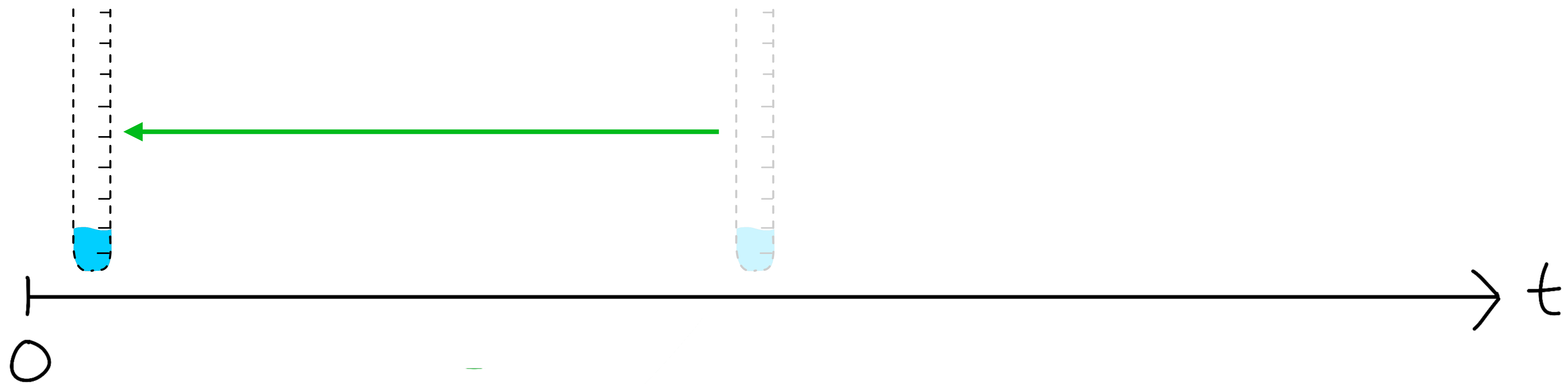defined by a *boost function $b(x) \geq 0$* that maps **age** to boost

**Which boost function is optimal?**

# The GittinsBoost policy

defined by a *boost function $b(x) \geq 0$* that maps **age** to boost

**Which boost function is optimal?**

choosing:

$$b(x) = \frac{1}{\gamma} \log \left( \sup_{y > x} \frac{\mathbf{E}[e^{\gamma S} \mathbf{1}(S \leq y) \mid S > x]}{\mathbf{E}[e^{\gamma((S \wedge y) - x)} - 1 \mid S > x]} \right)$$

gets us a strongly optimal policy in the class of policies that don't use job size information.

# The GittinsBoost policy

defined by a *boost function $b(x) \geq 0$* that maps **age** to boost

**Which boost function is optimal?**

choosing:

looks similar to the Gittins rank function…

$$b(x) = \frac{1}{\gamma} \log \left( \sup_{y > x} \frac{\mathbf{E}[e^{\gamma S} \mathbf{1}(S \leq y) \mid S > x]}{\mathbf{E}[e^{\gamma((S \wedge y) - x)} - 1 \mid S > x]} \right)$$

gets us a strongly optimal policy in the class of policies that don't use job size information.

# How might we discover an optimal policy?



random
arrivals

# How might we discover an optimal policy?



random arrivals

# How might we discover an optimal policy?



ignore future arrivals

random arrivals

# How might we discover an optimal policy?

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | | | |
| Batch Objective | | | |
| Optimal Policy | | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| **Queue Objective** | $\mathbf{E}_\pi[T]$ <br> **w/ known sizes** | | |
| **Batch Objective** | | | |
| **Optimal Policy** | | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$<br><br>**w/ known sizes** | | |
| Batch Objective | $\dfrac{1}{N}\displaystyle\sum_{i=1}^{N} T_i$<br><br>**w/ known sizes** | | |
| Optimal Policy | | | |

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$<br><br>**w/ known sizes** | | | |
| Batch Objective | $\sum_{i=1}^{N} T_i$<br><br>**w/ known sizes** | | | |
| Optimal Policy | | | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br> **w/ known sizes** | | |
| Batch Objective | $\displaystyle\sum_{i=1}^{N}(D_i - A_i)$ <br> **w/ known sizes** | | |
| Optimal Policy | | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br> **w/ known sizes** | | |
| Batch Objective | $\displaystyle\sum_{i=1}^{N} D_i - \sum_{i=1}^{N} A_i$ <br> **w/ known sizes** | | |
| Optimal Policy | | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br> **w/ known sizes** | | |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br> **w/ known sizes** | | |
| Optimal Policy | | | |

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br> **w/ known sizes** | | | |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br> **w/ known sizes** | | | |
| Optimal Policy | **SRPT** | | | |

# What is the optimal policy for the batch problem?

| | $\mathbf{E}_\pi[T]$ w/ known sizes | $\mathbf{E}_\pi[T]$ w/ unknown sizes | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br> w/ unknown sizes | | |
| Batch Objective | $\displaystyle\sum_{i=1}^{N} D_i$ <br> w/ known sizes | | | |
| Optimal Policy | **SRPT** | | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | |
| Batch Objective | $\displaystyle\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\displaystyle\mathbf{E}_\pi\!\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | |
| Optimal Policy | **SRPT** | | |

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | | |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | | |
| Optimal Policy | **SRPT** | **Gittins** | | |

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | $\lim_{t \to \infty} \mathbf{P}[T > t]$ <br><br> w/ known sizes | |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | | |
| Optimal Policy | **SRPT** | **Gittins** | | |

# What is the optimal policy for the batch problem?

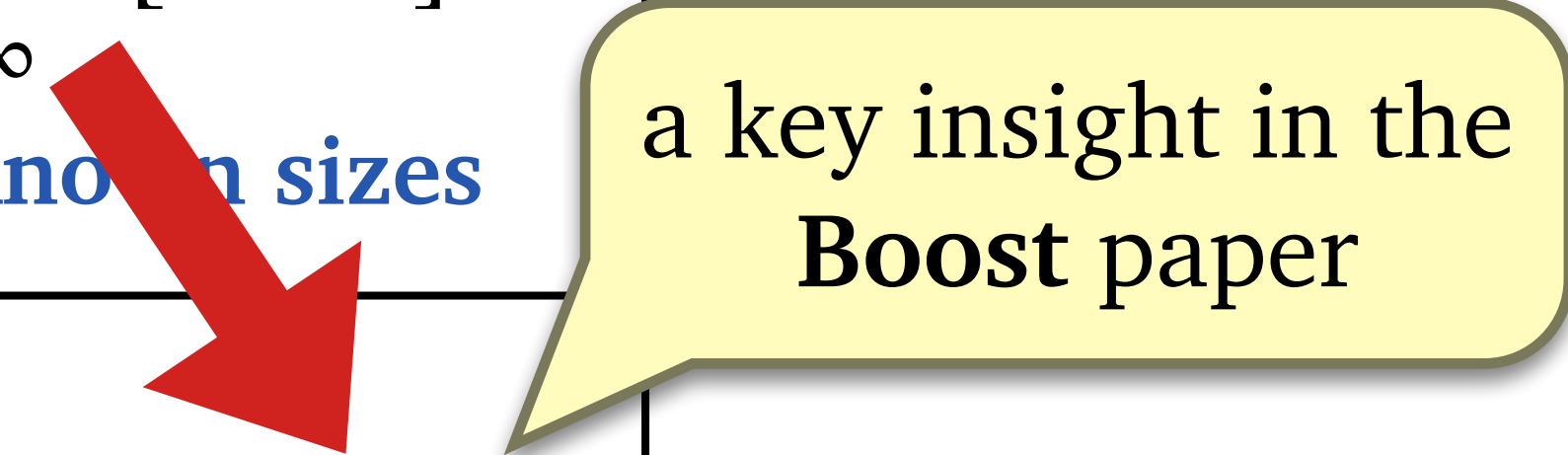| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ known sizes |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$ <br><br> w/ known sizes |
| Optimal Policy | **SRPT** | **Gittins** | |

a key insight in the **Boost** paper

# What is the optimal policy for the batch problem?

| | | | |
|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ known sizes |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$ <br><br> w/ known sizes |
| Optimal Policy | **SRPT** | **Gittins** | |

a key insight in the **Boost** paper

# What is the optimal policy for the batch problem?

| | $\mathbf{E}_\pi[T]$ <br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br> w/ known sizes | |
|---|---|---|---|---|
| **Queue Objective** | | | | |
| **Batch Objective** | $\sum_{i=1}^{N} D_i$ <br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br> w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i}\, e^{-\gamma A_i}$ <br> w/ known sizes | a key insight in the **Boost** paper |
| **Optimal Policy** | **SRPT** | **Gittins** | **Boost** | |

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| **Queue Objective** | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ known sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ unknown sizes |
| **Batch Objective** | $\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$ <br><br> w/ known sizes | |
| **Optimal Policy** | **SRPT** | **Gittins** | **Boost** | |

# What is the optimal policy for the batch problem?

| | $\mathbf{E}_\pi[T]$ <br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br> w/ known sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br> w/ unknown sizes |
|---|---|---|---|---|
| **Queue Objective** | | | | |
| **Batch Objective** | $\sum_{i=1}^{N} D_i$ <br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br> w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$ <br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}\right]$ <br> w/ unknown sizes |
| **Optimal Policy** | **SRPT** | **Gittins** | **Boost** | |

# What is the optimal policy for the batch problem?
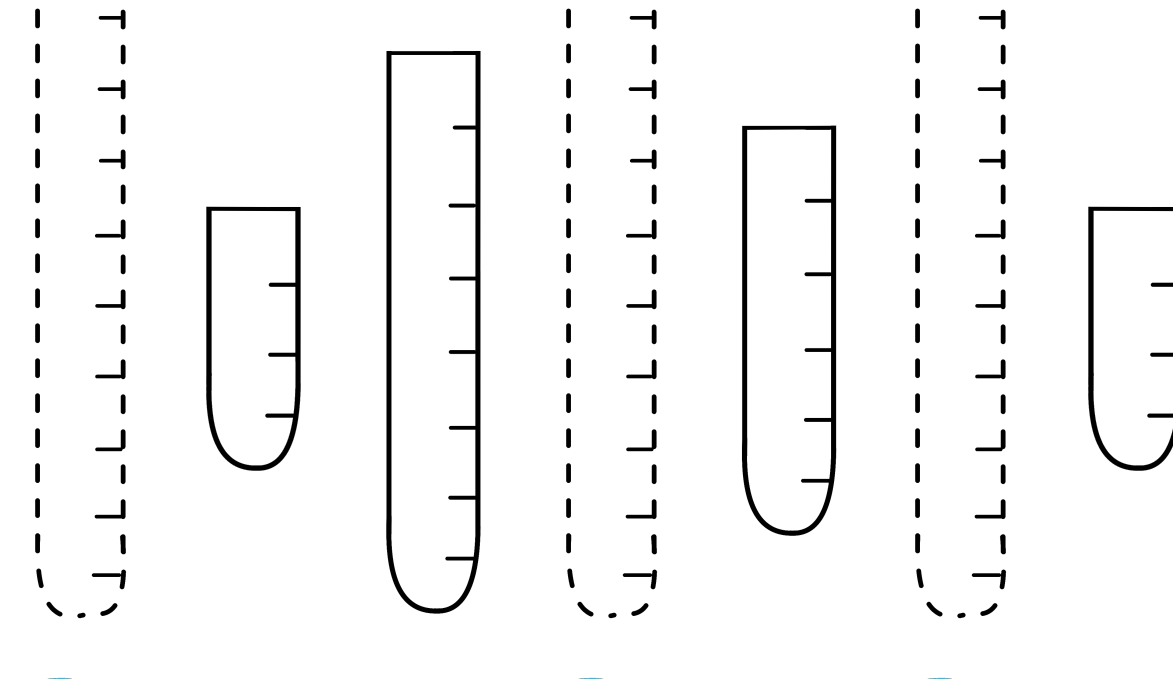
| | | | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$<br><br>w/ known sizes | $\mathbf{E}_\pi[T]$<br><br>w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$<br><br>w/ known sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$<br><br>w/ unknown sizes |
| Batch Objective | $\sum_{i=1}^{N} D_i$<br><br>w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$<br><br>w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$<br><br>w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}\right]$<br><br>w/ unknown sizes |
| Optimal Policy | **SRPT** | **Gittins** | **Boost** | **GittinsBoost** |

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ known sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ unknown sizes |
| Batch Objective | $\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i\right]$ <br><br> w/ unknown sizes | $\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$ <br><br> w/ known sizes | $\mathbf{E}_\pi\left[\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}\right]$ <br><br> w/ unknown sizes |
| Optimal Policy | **SRPT** | **Gittins** | **Boost** | **GittinsBoost** |

*All of these are in the Gittins family of policies!*

# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:
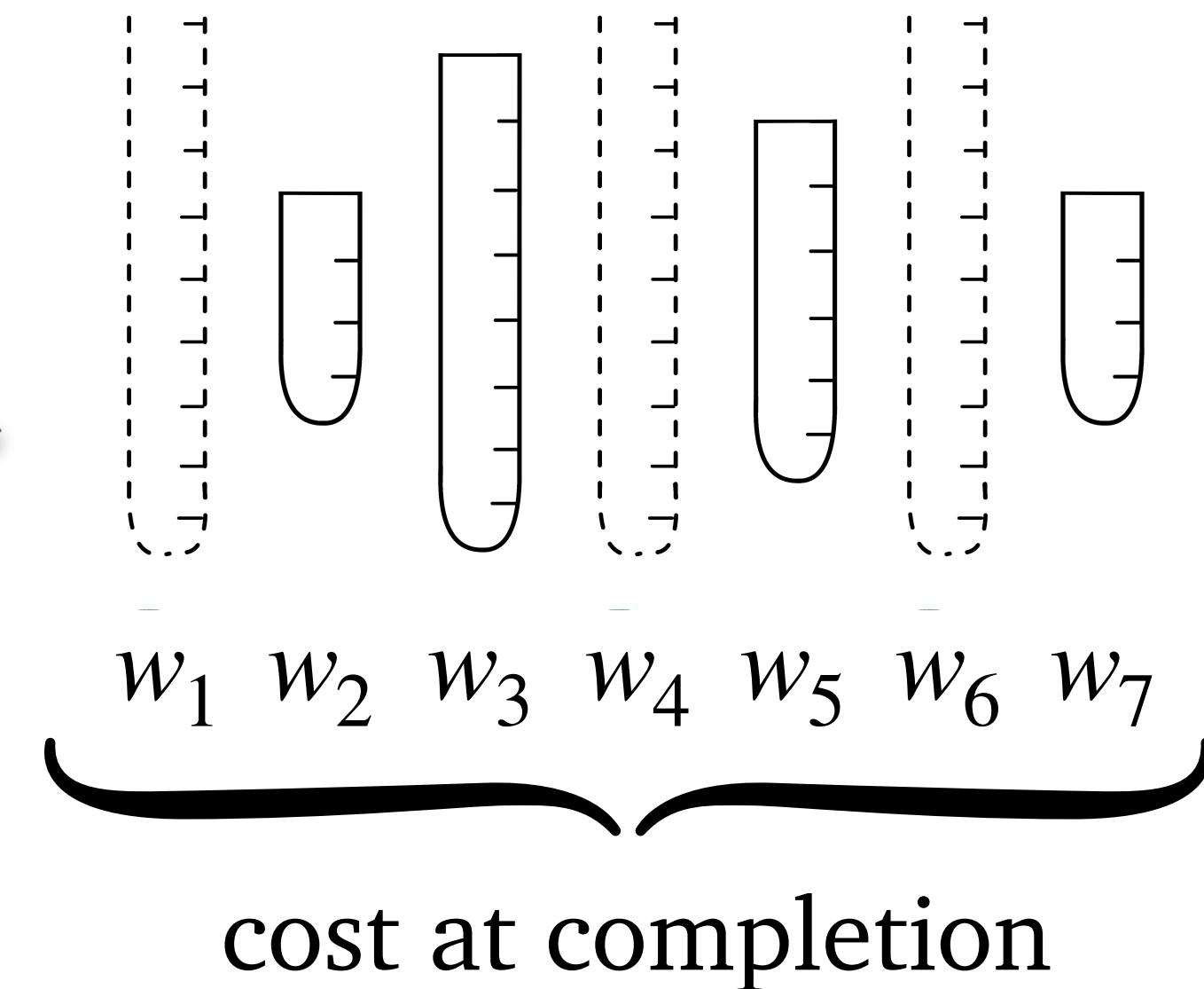
**job sizes independent**
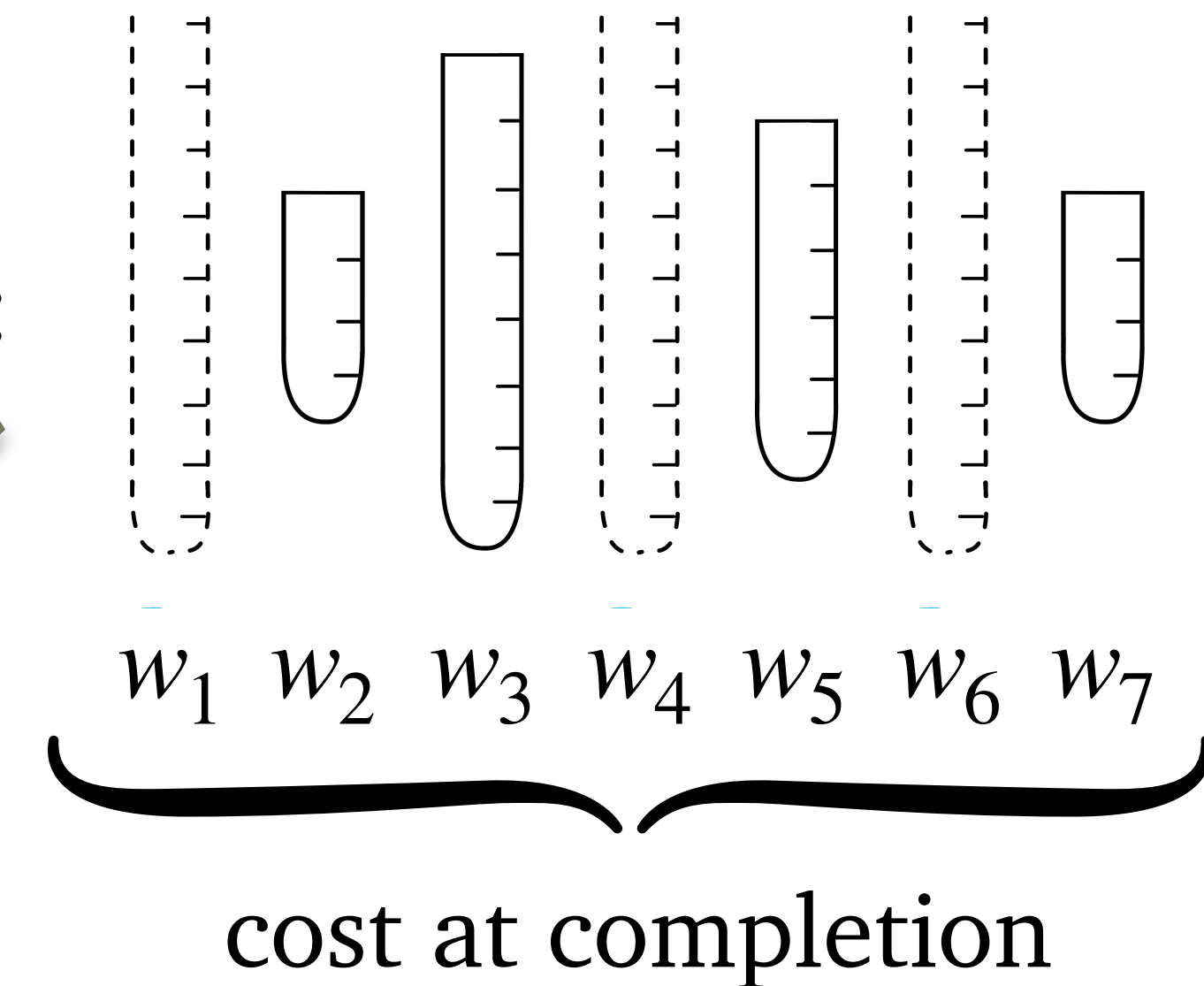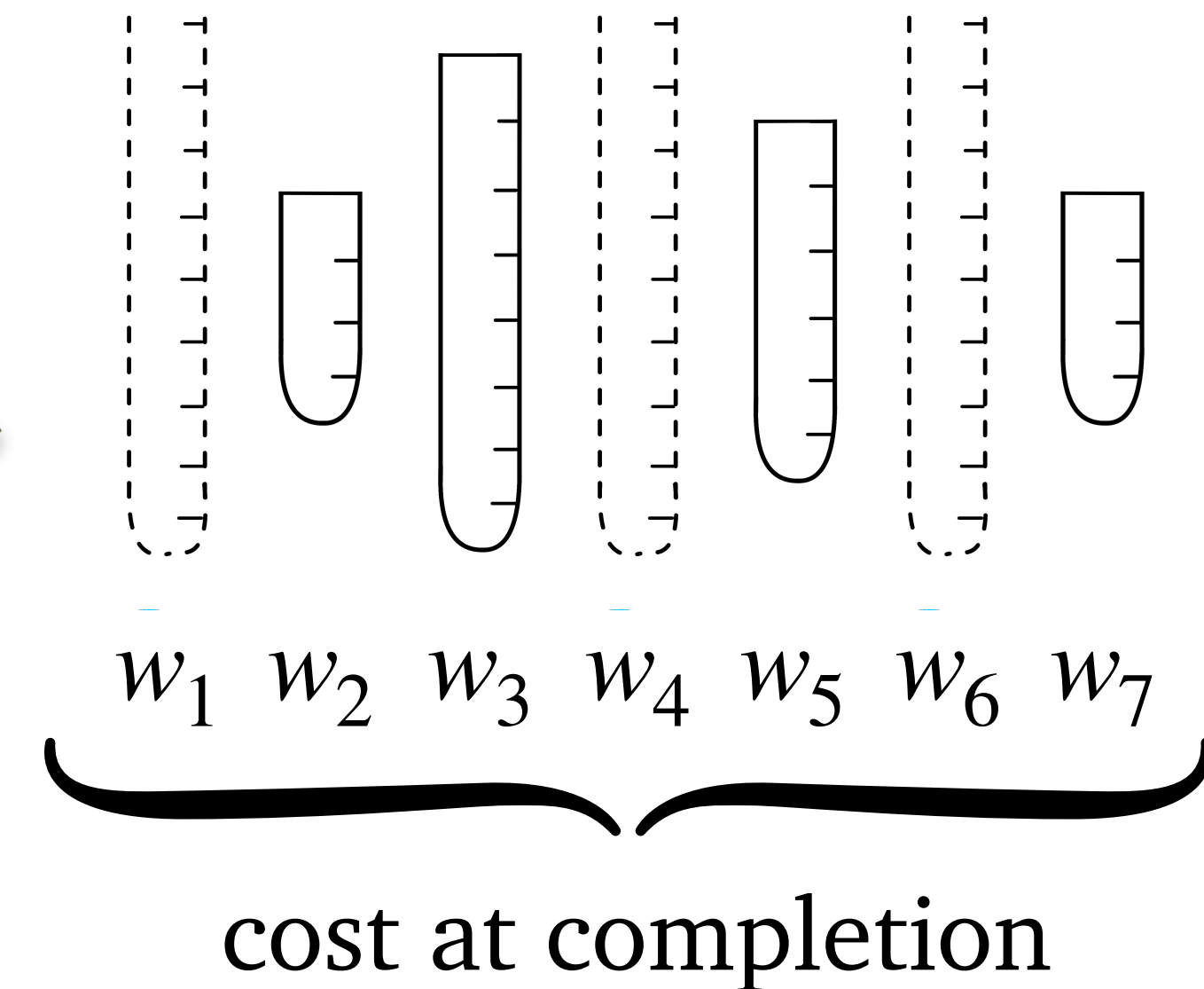
# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:

**job sizes independent**

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7$

cost at completion

# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:

job sizes independent

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$ $w_7$

cost at completion

with objective:

discounting

$$\text{minimize} \ \ \mathbf{E}_\pi\left[\sum_{i=1}^{N} e^{-\beta D_i} w_i\right]$$

$(\beta > 0)$

# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:

job sizes independent

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7$

cost at completion

with objective:

discounting

$$\text{minimize} \quad \mathbf{E}_\pi\left[ \sum_{i=1}^{N} e^{-\beta D_i} w_i \right] \quad \text{or} \quad \text{minimize} \quad \mathbf{E}_\pi\left[ \sum_{i=1}^{N} D_i w_i \right]$$

$$(\beta > 0) \qquad\qquad\qquad\qquad\qquad\qquad \text{``}(\beta = 0)\text{''}$$

# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:

**job sizes independent**

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7$

cost at completion

with objective:

discounting

$$\text{minimize } \mathbf{E}_\pi \left[ \sum_{i=1}^{N} e^{-\beta D_i} w_i \right] \quad \text{or} \quad \text{minimize } \mathbf{E}_\pi \left[ \sum_{i=1}^{N} D_i w_i \right] \quad \text{or} \quad \text{minimize } \mathbf{E}_\pi \left[ \sum_{i=1}^{N} e^{-\beta D_i} w_i \right]$$

$$(\beta > 0) \qquad\qquad\qquad \text{``}(\beta = 0)\text{''} \qquad\qquad \textbf{new-ish!} \quad (\beta < 0)$$

# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:



**job sizes independent**

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7$$

cost at completion

with objective:

**discounting**

**negative discounting = inflation**

$$\text{minimize } \mathbf{E}_\pi\left[\sum_{i=1}^{N} e^{-\beta D_i} w_i\right] \quad \text{or} \quad \text{minimize } \mathbf{E}_\pi\left[\sum_{i=1}^{N} D_i w_i\right] \quad \text{or} \quad \text{minimize } \mathbf{E}_\pi\left[\sum_{i=1}^{N} e^{-\beta D_i} w_i\right]$$

$$(\beta > 0) \qquad\qquad\qquad\qquad\qquad \text{``}(\beta = 0)\text{''} \qquad\qquad \textbf{new-ish!} \qquad (\beta < 0)$$

# What is the Gittins family of policies?

Gittins policies solve the family of batch problems:

**job sizes independent**

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7$

cost at completion

with objective:

discounting

$$\text{minimize } \mathbf{E}_\pi\left[\sum_{i=1}^N e^{-\beta D_i} w_i\right] \quad \text{or} \quad \text{minimize } \mathbf{E}_\pi\left[\sum_{i=1}^N D_i w_i\right] \quad \text{or} \quad \text{minimize } \mathbf{E}_\pi\left[\sum_{i=1}^N e^{-\beta D_i} w_i\right]$$

$e^{\gamma D_i} e^{-\gamma A_i}$

$(\beta > 0)$ $\text{``}(\beta = 0)\text{''}$ **new-ish!** $(\beta < 0)$

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| **Queue Objective** | $\mathbf{E}_\pi[T]$ <br><br> w/ known sizes | $\mathbf{E}_\pi[T]$ <br><br> w/ unknown sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ known sizes | $\lim_{t\to\infty} \mathbf{P}[T > t]$ <br><br> w/ unknown sizes |
| **Batch Objective** | $\displaystyle\sum_{i=1}^{N} D_i$ <br><br> w/ known sizes | $\displaystyle\mathbf{E}_\pi\Big[ \sum_{i=1}^{N} D_i \Big]$ <br><br> w/ unknown sizes | $\displaystyle\sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$ <br><br> w/ known sizes | $\displaystyle\mathbf{E}_\pi\Big[ \sum_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i} \Big]$ <br><br> w/ unknown sizes |
| **Optimal Policy** | **SRPT** | **Gittins** | **Boost** | **GittinsBoost** |

*All of these are in the Gittins family of policies!*

# What is the optimal policy for the batch problem?

| | | | | |
|---|---|---|---|---|
| Queue Objective | $\mathbf{E}_\pi[T]$<br><br>w/ known sizes | $\mathbf{E}_\pi[T]$<br><br>w/ unknown sizes | $\lim\limits_{t\to\infty} \mathbf{P}[T > t]$<br><br>w/ known sizes | $\lim\limits_{t\to\infty} \mathbf{P}[T > t]$<br><br>w/ unknown sizes |
| Batch Objective | $\sum\limits_{i=1}^{N} D_i$<br><br>w/ known sizes | $\mathbf{E}_\pi\left[\sum\limits_{i=1}^{N} D_i\right]$<br><br>w/ unknown sizes | $\sum\limits_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}$<br><br>w/ known sizes | $\mathbf{E}_\pi\left[\sum\limits_{i=1}^{N} e^{\gamma D_i} e^{-\gamma A_i}\right]$<br><br>w/ unknown sizes |
| Optimal Policy | **SRPT** | **Gittins** | **Boost** | **GittinsBoost** |

*All of these are in the Gittins family of policies!*

**How do we show optimality in the queue setting?**

# Boost optimality in the queue setting

server idle

server idle

busy period

# Boost optimality in the queue setting



server idle

server idle

$O$

$t$

busy period

batch problem:

# GittinsBoost optimality in the queue setting

# GittinsBoost optimality in the queue setting



server idle

server idle

$O$

$t$

137

# **GittinsBoost** optimality in the queue setting

server idle

server idle

batch problem:

# GittinsBoost optimality in the queue setting



server idle

server idle

job sizes are now correlated

batch problem:

# **GittinsBoost optimality in the queue setting**

server idle                    server idle



job sizes are now correlated

## **Example**

$\lambda = \varepsilon \ll 1$ and $S = \mathrm{Unif}\{1, \varepsilon\}$

batch problem:

# GittinsBoost optimality in the queue setting



**Example**

$\lambda = \varepsilon \ll 1$ and $S = \mathrm{Unif}\{1, \varepsilon\}$

3 jobs: $A_1 = 0$, $A_2 = \varepsilon^2$, $A_3 = 1$

server idle

server idle

job sizes are now correlated

batch problem:

# GittinsBoost optimality in the queue setting



server idle

server idle

job sizes are now correlated

**Example**

$\lambda = \varepsilon \ll 1$ and $S = \mathrm{Unif}\{1, \varepsilon\}$

3 jobs: $A_1 = 0$, $A_2 = \varepsilon^2$, $A_3 = 1$

If $S_1 = \varepsilon$ then $S_2 = 1$

batch problem:

# GittinsBoost optimality in the queue setting



server idle

server idle

job sizes are now correlated

Gittins policy is not optimal for
this correlated batch problem

batch problem:

# **GittinsBoost** optimality in the queue setting



server idle

server idle

job sizes are now correlated

$O$

$t$

**Gittins policy is not optimal for this correlated batch problem**

batch problem:

*main technical challenge*: showing optimality in queue setting

# What was our approach?

# What was our approach?

**Boosted Arrival Time**

non-preemptible    non-preemptible

**age**

**Boosted Arrival Time**

**age**

# What was our approach?



Boosted Arrival Time — non-preemptible — non-preemptible — age

Boosted Arrival Time — age

Boosted Arrival Time — age

# What was our approach?



**Boosted Arrival Time** (×3)

non-preemptible   non-preemptible

age   age   age

**Batch Setting Optimality:** all three policies are the same

# What was our approach?

**Boosted Arrival Time**

non-preemptible    non-preemptible

**age**

**Boosted Arrival Time**

**age**

**Boosted Arrival Time**

**age**

**Batch Setting Optimality:** all three policies are the same

**Queue Setting Optimality:** all three policies have the same asymptotic tail behavior

# Summary

# Summary

**Problem**



Schedule for $\mathbf{P}[T > t]$ as $t \to \infty$

# Summary



## Problem

Schedule for $\mathbf{P}[T > t]$ as $t \to \infty$
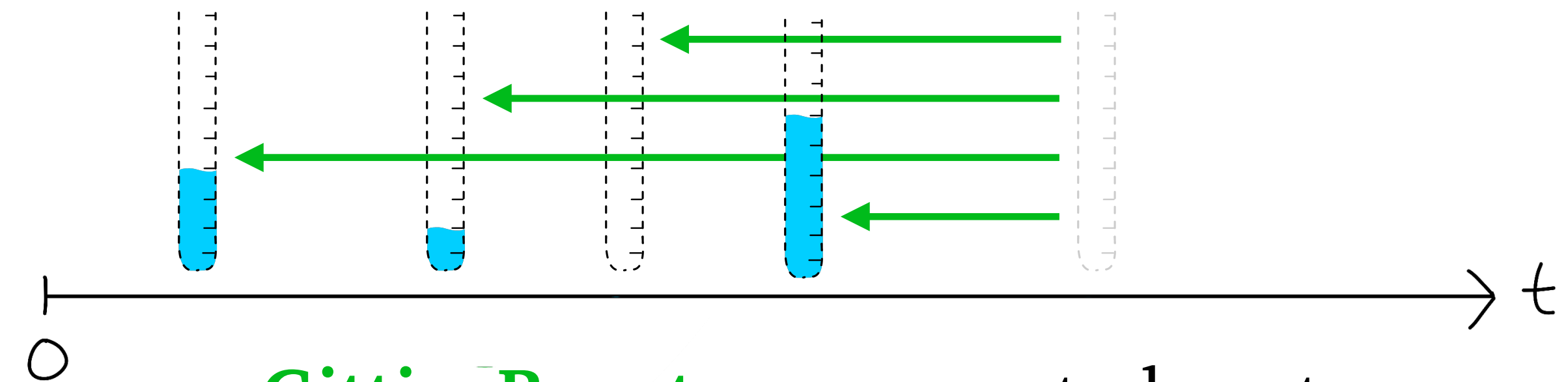
## Contribution

**GittinsBoost:** map **age** to boost

# Summary

## Problem



Schedule for $\mathbf{P}[T > t]$ as $t \to \infty$

## Contribution



**GittinsBoost:** map **age** to boost

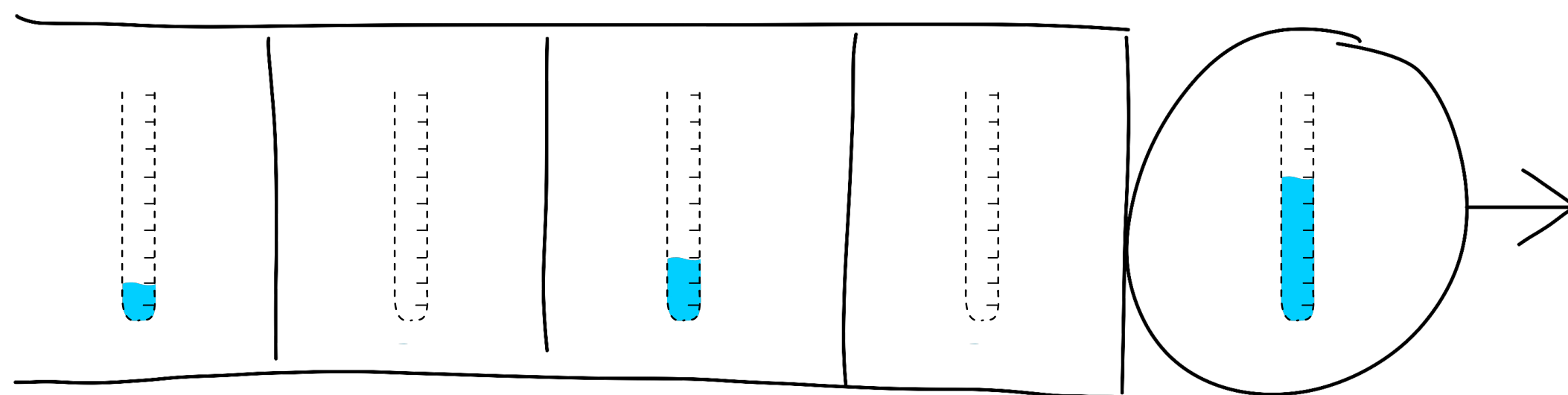## Main Ideas



batch problem:

# Summary

## Problem



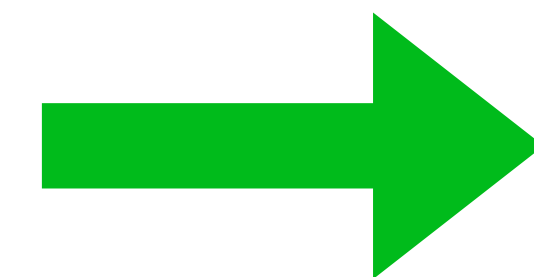Schedule for $\mathbf{P}[T > t]$ as $t \to \infty$

## Contribution



**GittinsBoost:** map **age** to boost

## Main Ideas
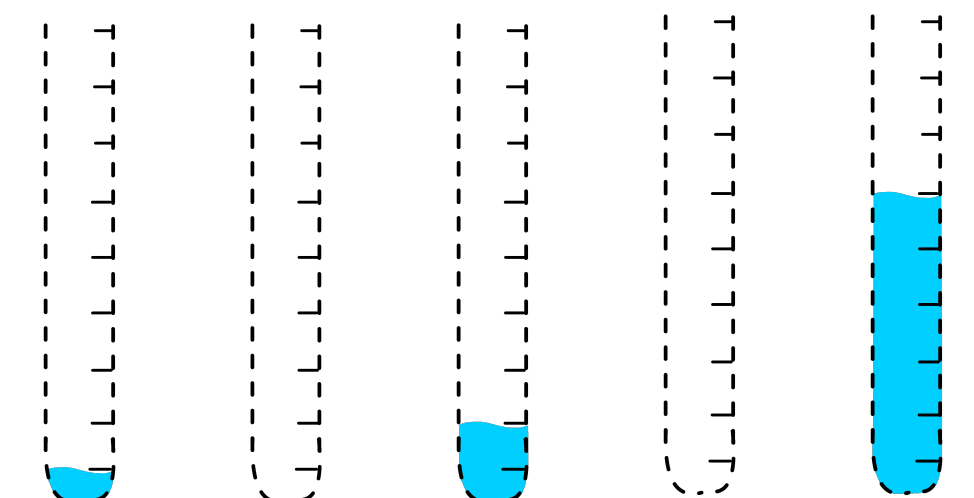


batch problem:

queue optimality ⬅ batch optimality

*main technical challenge*