

# *A* **New Toolbox** *for* **Scheduling Theory**

Ziv Scully

*Harvard & MIT → Cornell*

`zivscully@cornell.edu`

`https://ziv.codes`



# Collaborators



Mor Harchol-Balter  
*CMU, thesis advisor*



Guy Blelloch  
*CMU, thesis advisor*



Alan Scheller-Wolf  
*CMU*



Isaac Grosf  
*CMU*

Adam Wierman  
*Caltech*

Onno Boxma  
*TU/e (Eindhoven)*

Jan-Pieter Dorsman  
*UvA (Amsterdam)*

Lucas van Kreveld  
*UvA (Amsterdam)*

Michael Mitzenmacher  
*Harvard*

Sid Banerjee  
*Cornell*

Anupam Gupta  
*CMU*

Sahil Singla  
*Georgia Tech*

Haotian Jiang  
*University of Washington*

Kunhe Yang  
*UC Berkeley*

# Performance

# Performance

*efficiency throughput delay reliability safety value*

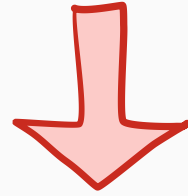


# Performance

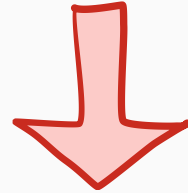
*efficiency*   *throughput*   *delay*   *reliability*   *safety*   *value*

this talk

# Contention

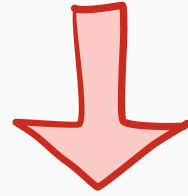


# Queueing

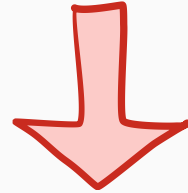


# Delay

# Contention



# Queueing



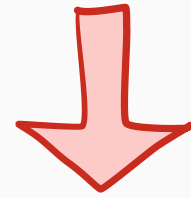
# Delay

healthcare

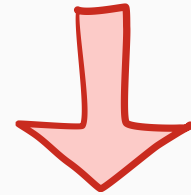
supply chains

**Contention**

healthcare



**Queueing**



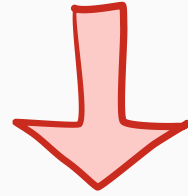
**Delay**

your local supermarket

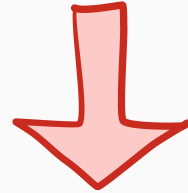
supply chains

**Contention**

healthcare



**Queueing**



**Delay**



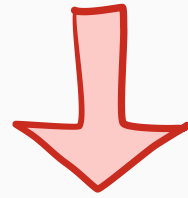
your local supermarket

supply chains

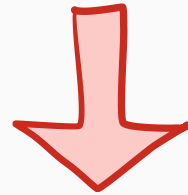
**Contention**

call centers

healthcare



**Queueing**



**Delay**

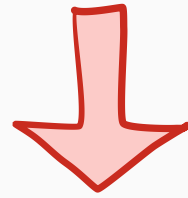
your local supermarket

supply chains

**Contention**

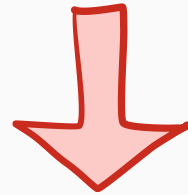
call centers

healthcare



**Queueing**

transportation



**Delay**

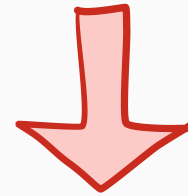
your local supermarket

supply chains

**Contention**

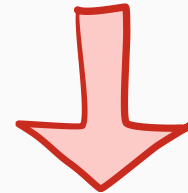
call centers

healthcare



**Queueing**

transportation



**Delay**

databases

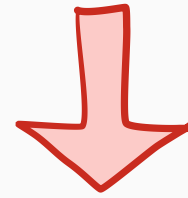
your local supermarket

supply chains

**Contention**

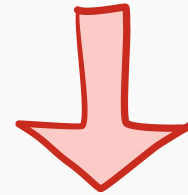
call centers

healthcare



**Queueing**

transportation



databases

**Delay**

networks

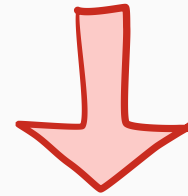
your local supermarket

supply chains

**Contention**

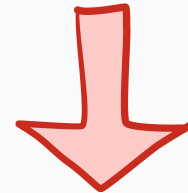
call centers

healthcare



**Queueing**

transportation



databases

**Delay**

networks

operating systems



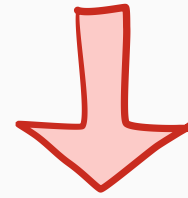
your local supermarket

supply chains

**Contention**

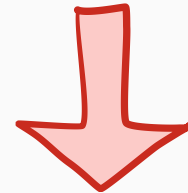
call centers

healthcare



**Queueing**

transportation



databases

computer architecture

**Delay**

networks

operating systems

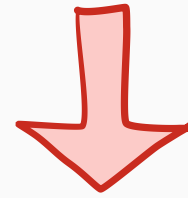
your local supermarket

supply chains

**Contention**

call centers

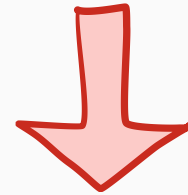
healthcare



**Queueing**

transportation

supercomputing



databases

computer architecture

**Delay**

networks

operating systems

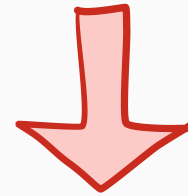
your local supermarket

supply chains

**Contention**

call centers

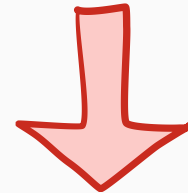
healthcare



**Queueing**

transportation

supercomputing



databases

computer architecture

**Delay**

networks

operating systems



*How to reduce delays?*

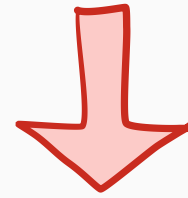
your local supermarket

supply chains

**Contention**

call centers

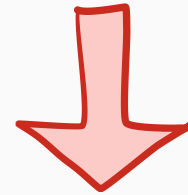
healthcare



**Queueing**

transportation

supercomputing



databases

computer architecture

**Delay**

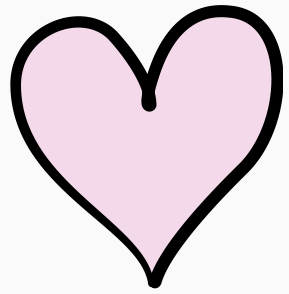
networks

operating systems



*How to reduce delays?*

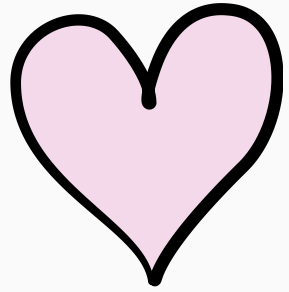
**Scheduling**



*Good news:*

scheduling can reduce delay





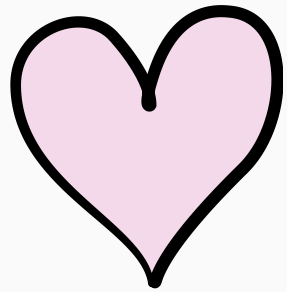
*Good news:*

scheduling can reduce delay



*Bad news:*

limited understanding of scheduling



*Good news:*

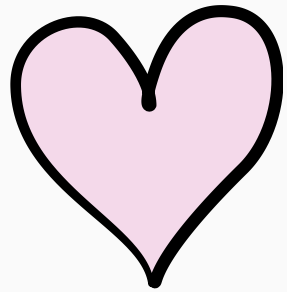
scheduling can reduce delay



*Bad news:*

limited understanding of scheduling

evaluation



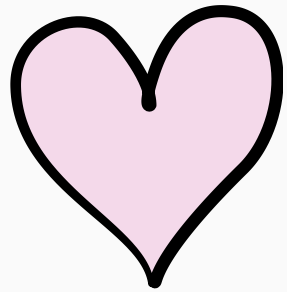
*Good news:*  
scheduling can reduce delay



*Bad news:*  
limited understanding of scheduling

design

evaluation



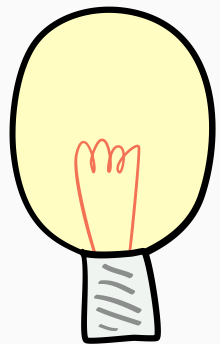
*Good news:*  
scheduling can reduce delay



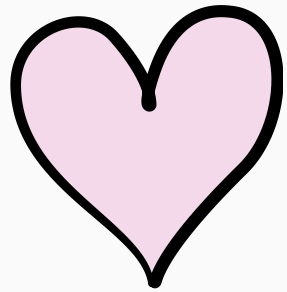
*Bad news:*  
limited understanding of scheduling

design

evaluation



*We need:*  
rigorous theory of scheduling



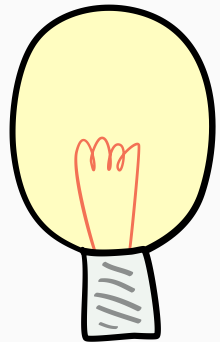
*Good news:*  
scheduling can reduce delay



*Bad news:*  
limited understanding of scheduling

design

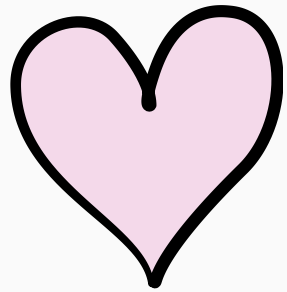
evaluation



*We need:*  
rigorous theory of scheduling

CS, EE, OR, OM





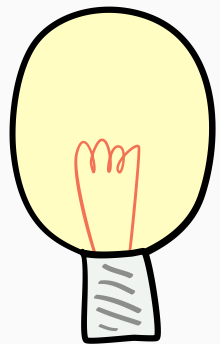
*Good news:*  
scheduling can reduce delay



*Bad news:*  
limited understanding of scheduling

design

evaluation



*We need:*  
rigorous theory of scheduling

applied math, probability,  
control, learning theory

CS, EE, OR, OM

# Theory lags behind practice

# Theory lags behind practice

**Uncertainty:**

unknown job sizes, noisy estimates

# Theory lags behind practice

## **Uncertainty:**

unknown job sizes, noisy estimates

## **Multiple servers:**

even FCFS hard to understand

# Theory lags behind practice

## **Uncertainty:**

unknown job sizes, noisy estimates

## **Multiple servers:**

even FCFS hard to understand

## **Preemption practicalities:**

limitations, overhead

# Theory lags behind practice



**This talk:**  
*new theoretical tools*

## **Uncertainty:**

unknown job sizes, noisy estimates

## **Multiple servers:**

even FCFS hard to understand

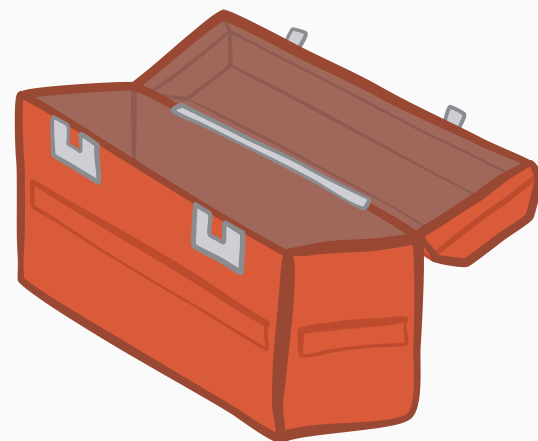
## **Preemption practicalities:**

limitations, overhead

# Outline: two new tools



# Outline: two new tools

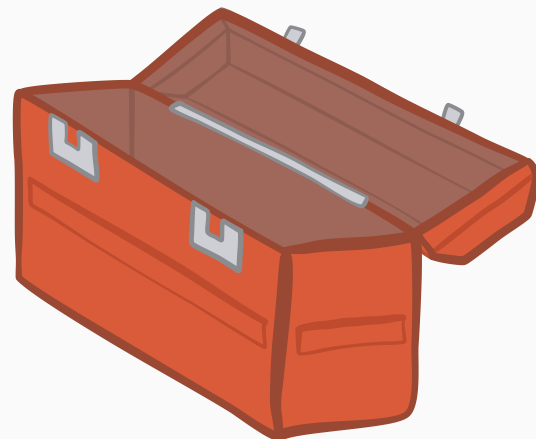




# Outline: two new tools



new unifying theory of  
single-server scheduling

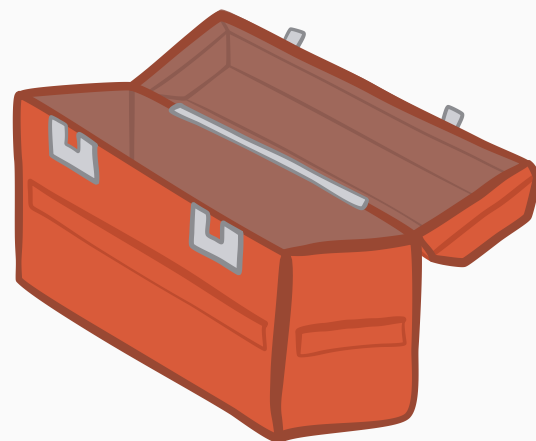


# Outline: two new tools



new unifying theory of  
single-server scheduling

greatly increases **number  
of policies** we can analyze



# Outline: two new tools

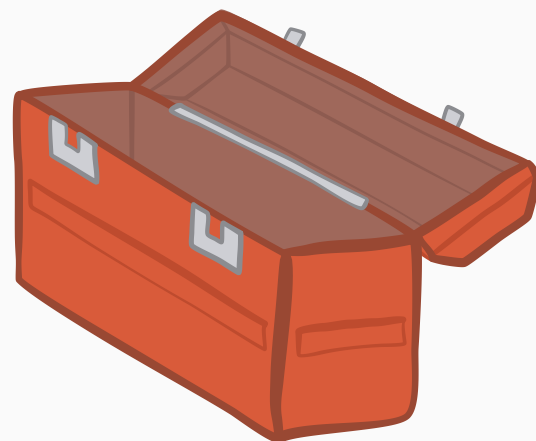


new unifying theory of  
single-server scheduling

greatly increases **number  
of policies** we can analyze

→ unknown sizes

→ preemption limitations

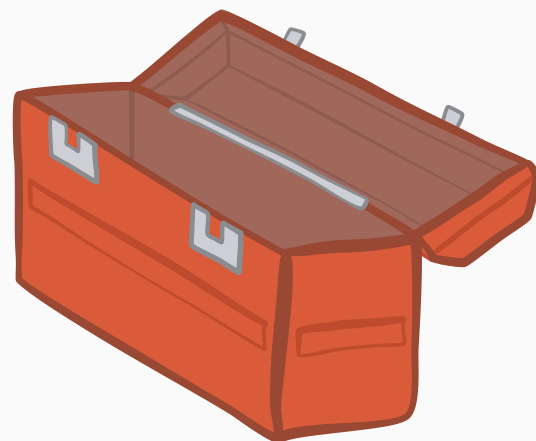


# Outline: two new tools



new unifying theory of  
single-server scheduling

- unknown sizes
- preemption limitations



# Outline: two new tools



## SOAP

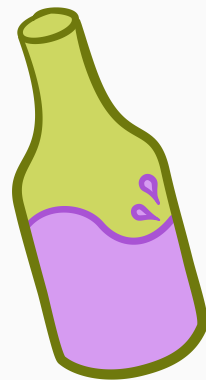
new unifying theory of  
single-server scheduling



unknown sizes

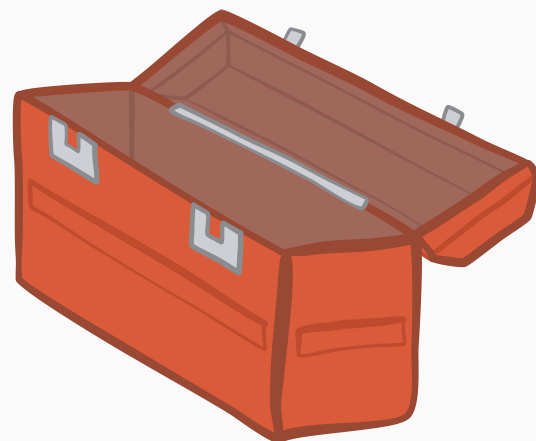


preemption limitations

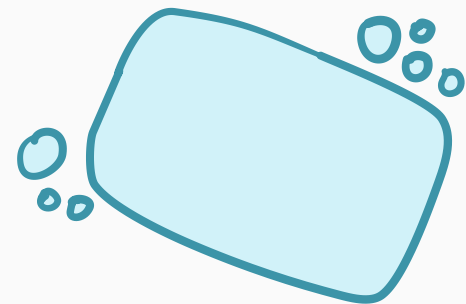


## WINE

new queueing identity to  
complement Little's Law



# Outline: two new tools



## SOAP

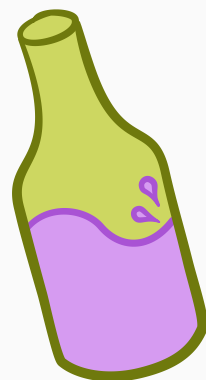
new unifying theory of single-server scheduling



unknown sizes



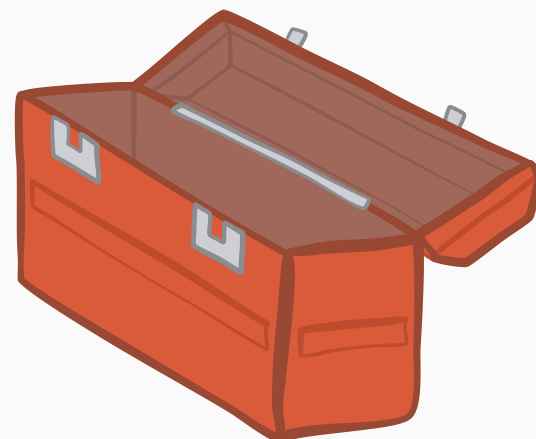
preemption limitations



## WINE

new queueing identity to complement Little's Law

reduces **complex systems** to related **simple systems**



# Outline: two new tools



new unifying theory of single-server scheduling

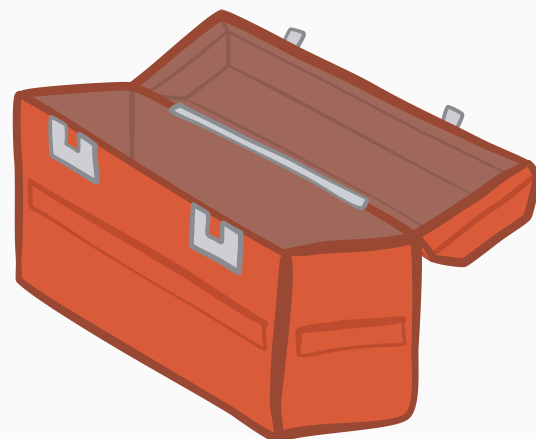
- unknown sizes
- preemption limitations



new queueing identity to complement Little's Law

reduces **complex systems** to related **simple systems**

- multiserver systems
- noisy size estimates



# Outline: two new tools



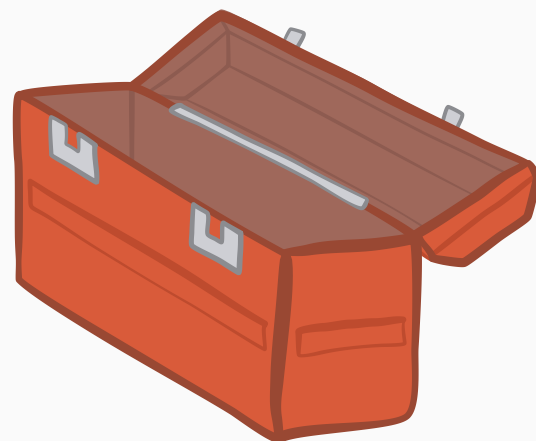
new unifying theory of  
single-server scheduling

- unknown sizes
- preemption limitations



new queueing identity to  
complement Little's Law

- multiserver systems
- noisy size estimates

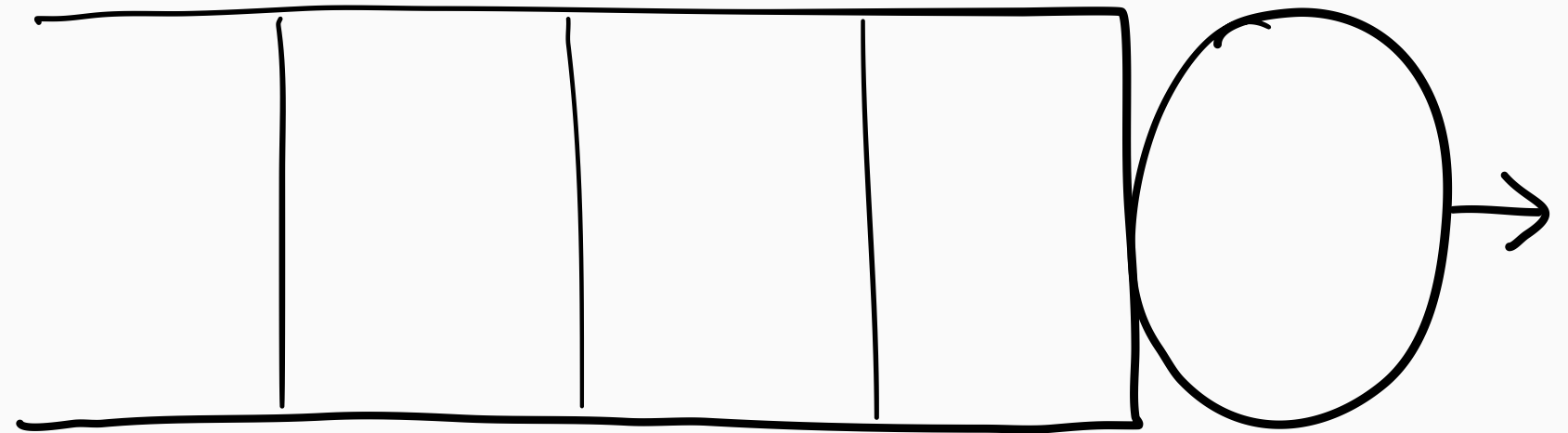






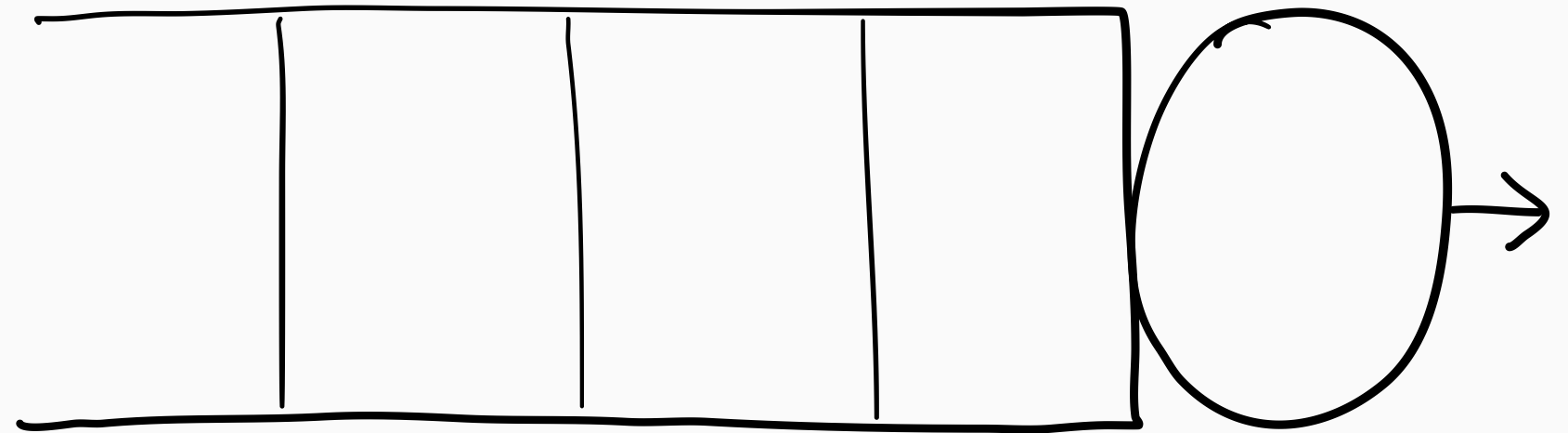
*Part 1*  
**SOAP**

# Single-server queueing system



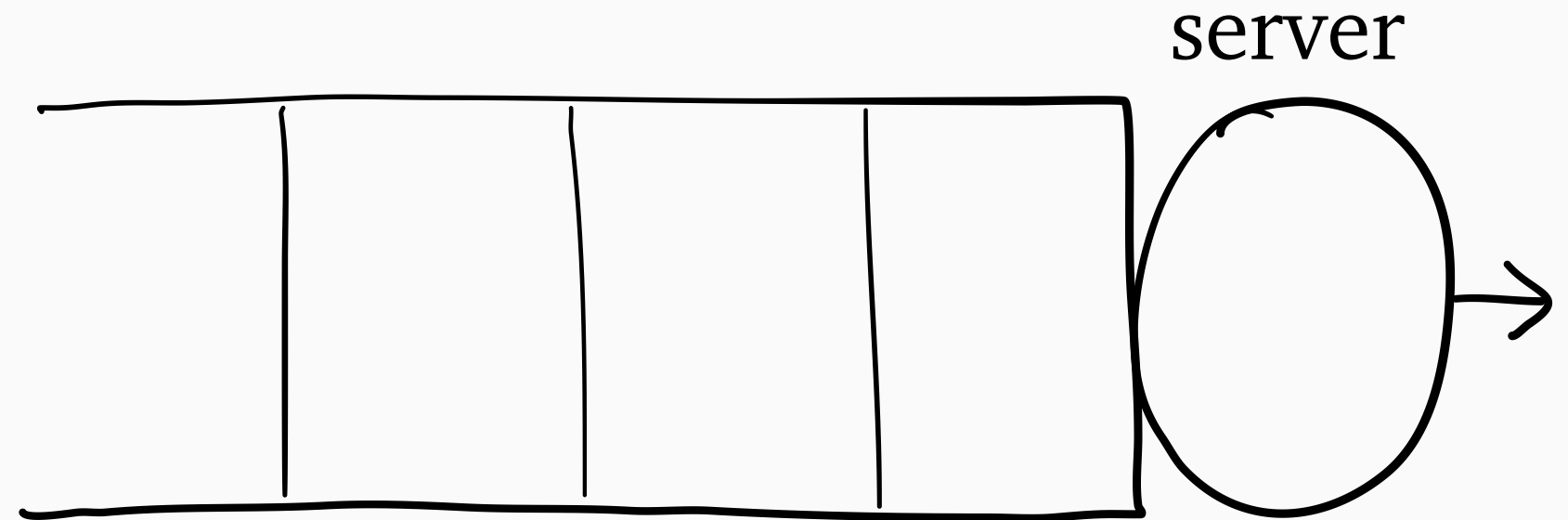
# Single-server queueing system

M/G/1



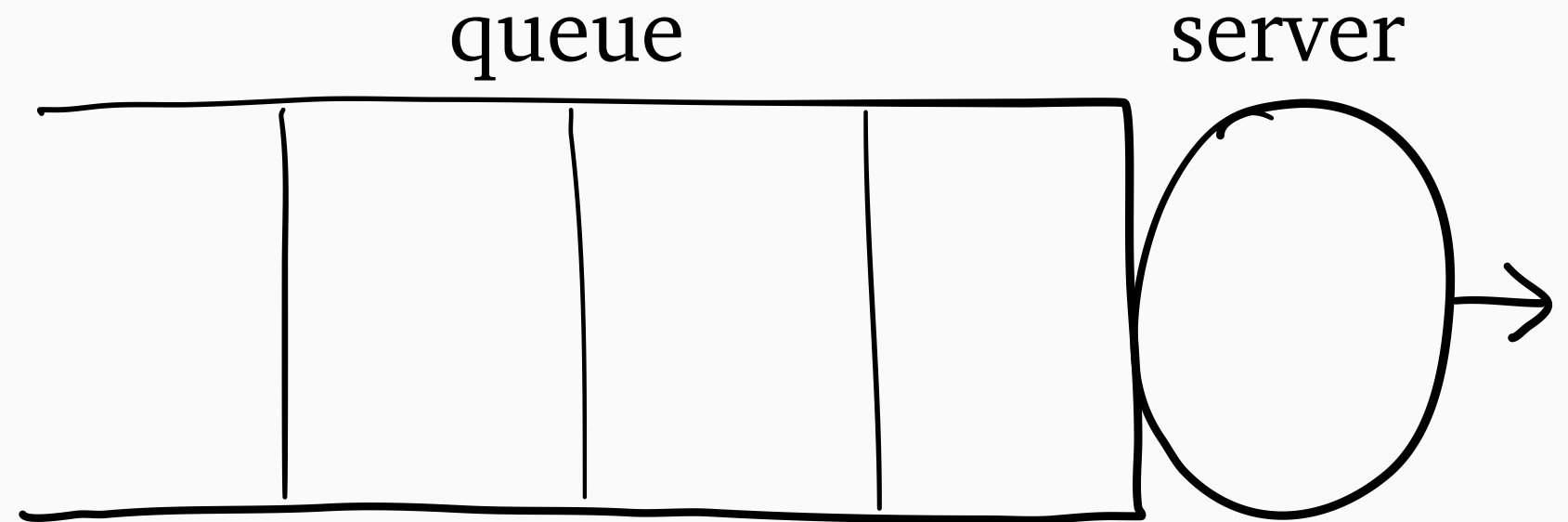
# Single-server queueing system

M/G/1



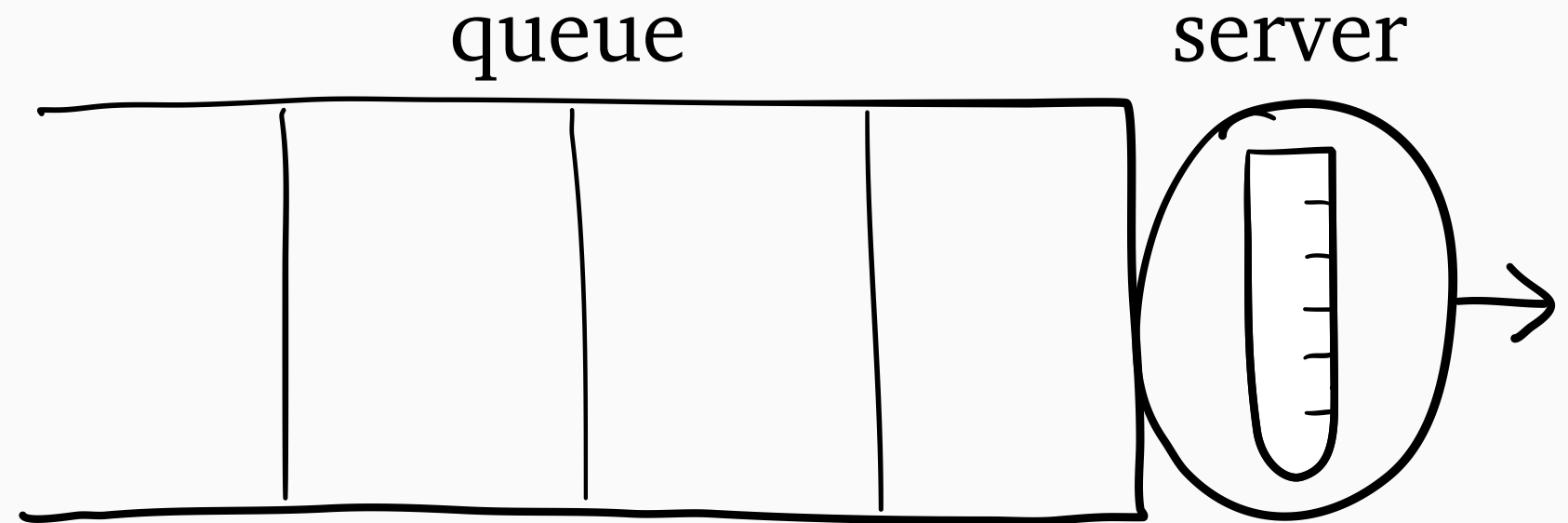
# Single-server queueing system

M/G/1

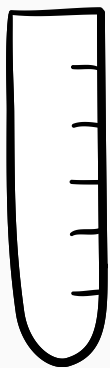


# Single-server queueing system

M/G/1

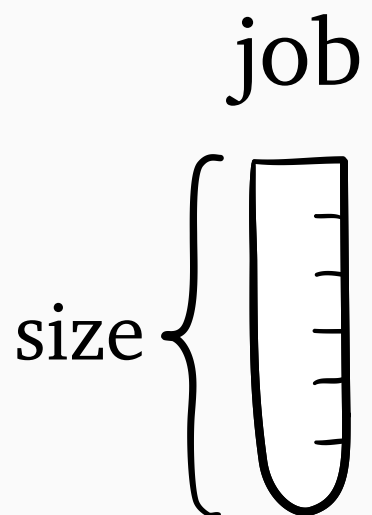
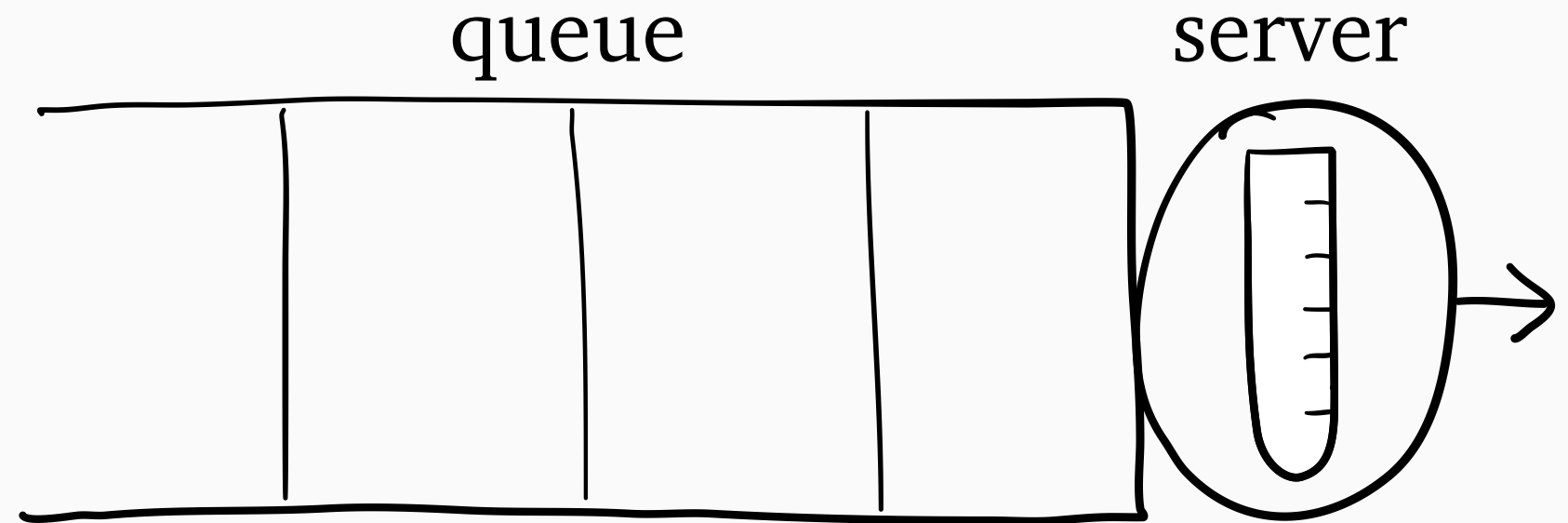


job



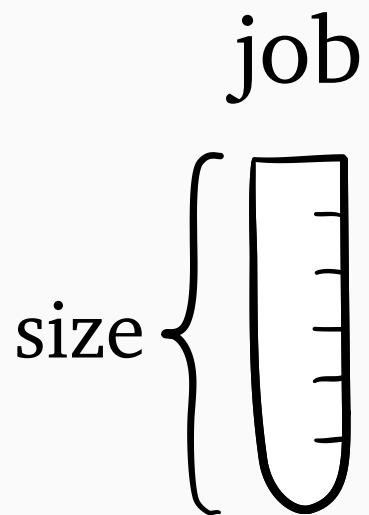
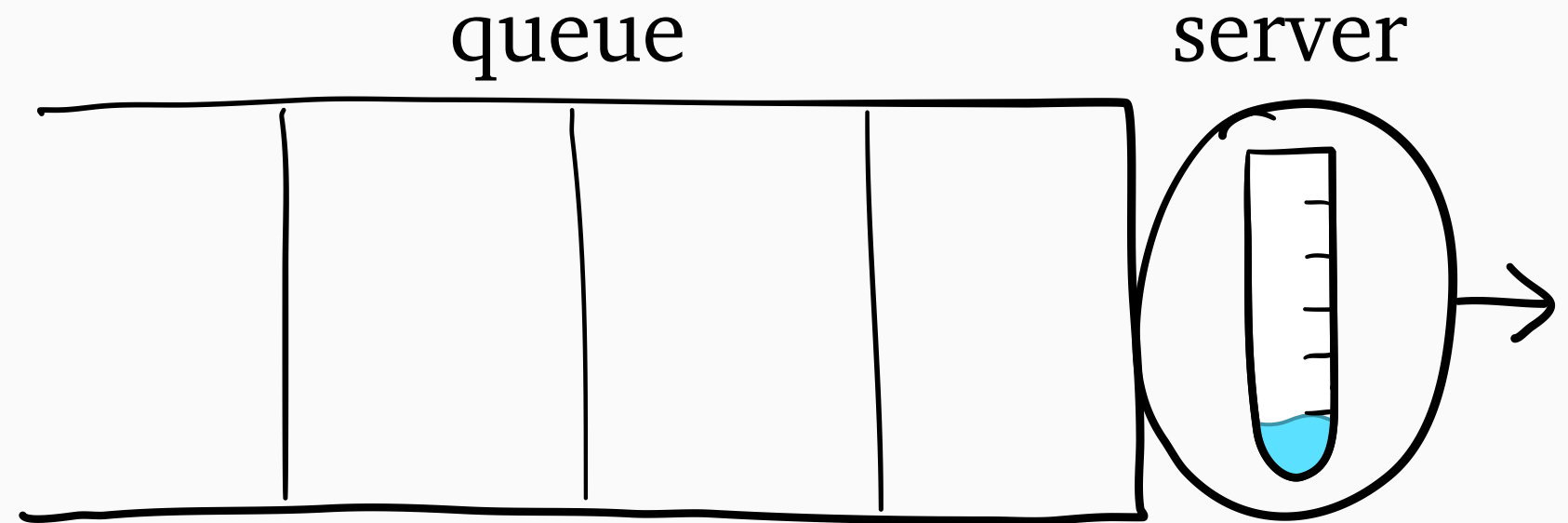
# Single-server queueing system

M/G/1



# Single-server queueing system

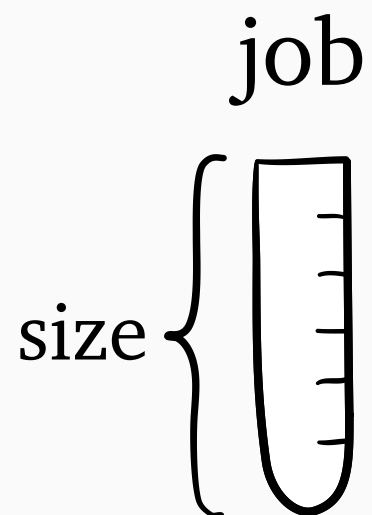
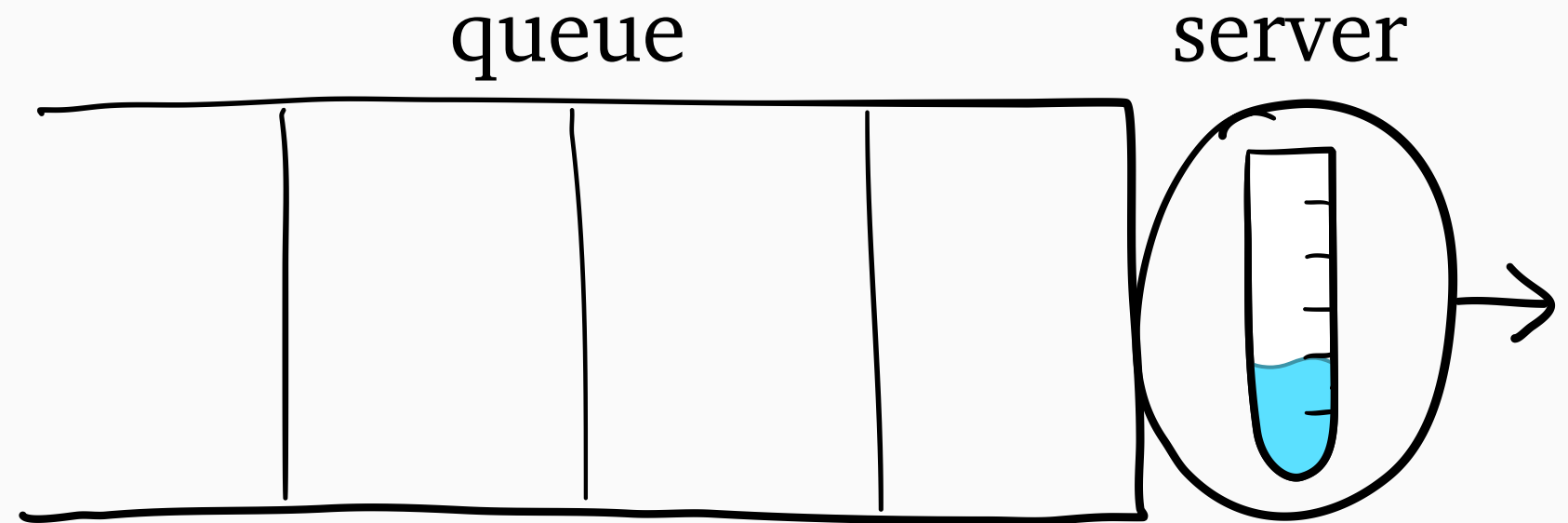
M/G/1





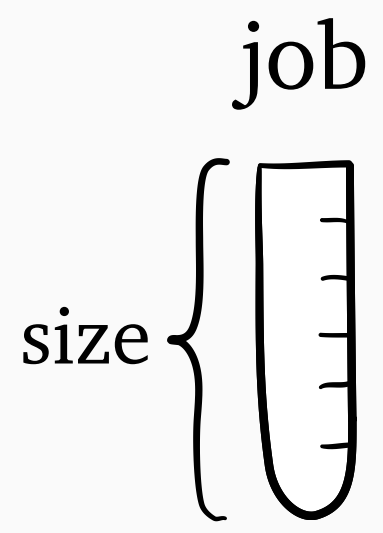
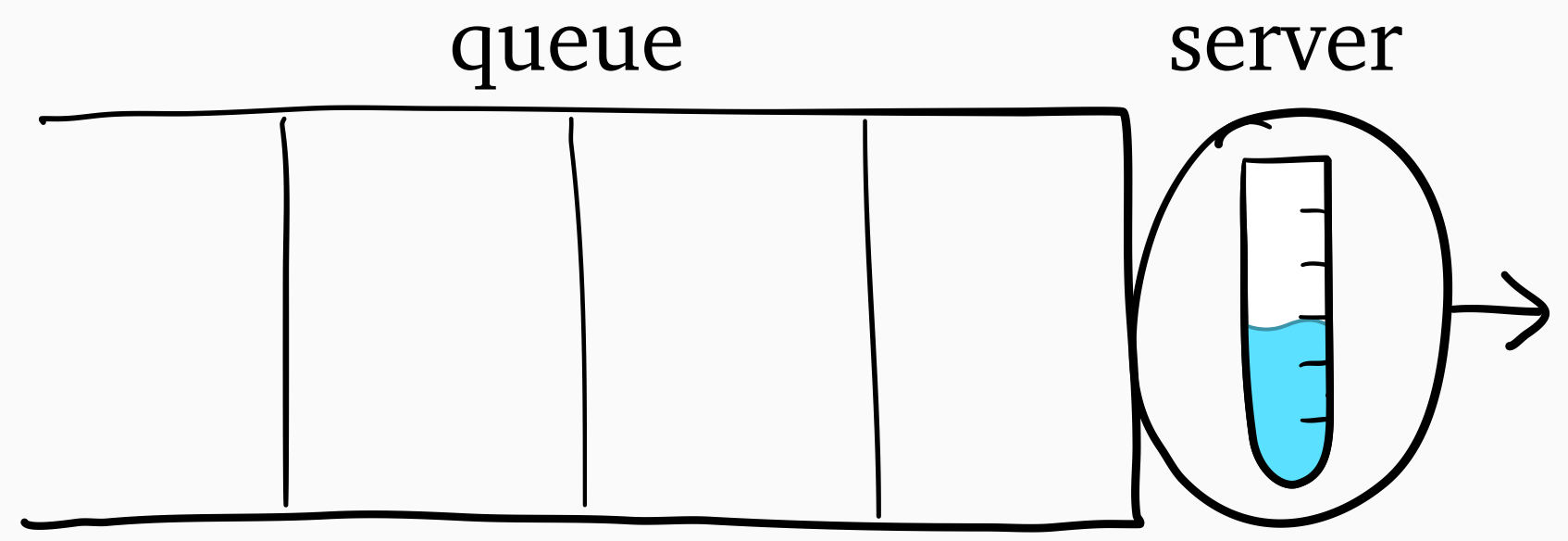
# Single-server queueing system

M/G/1



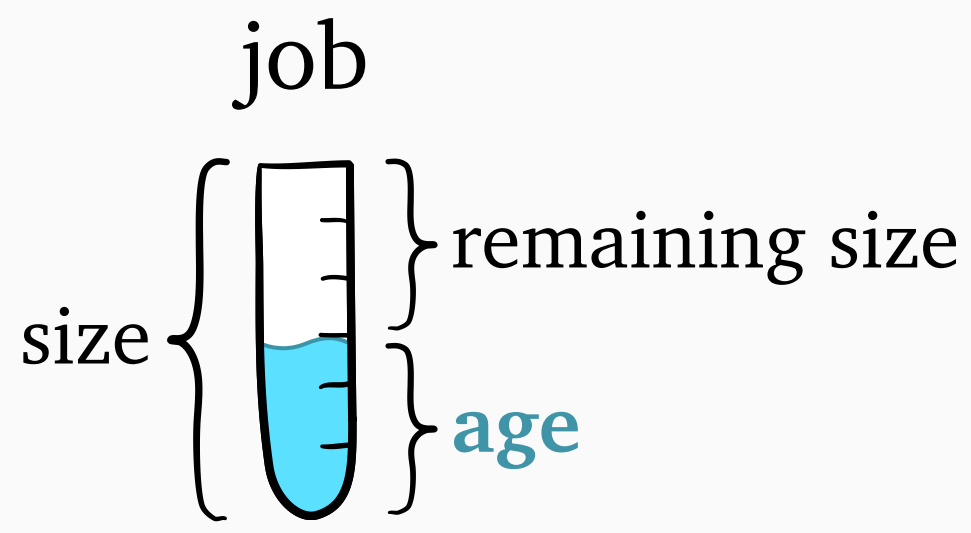
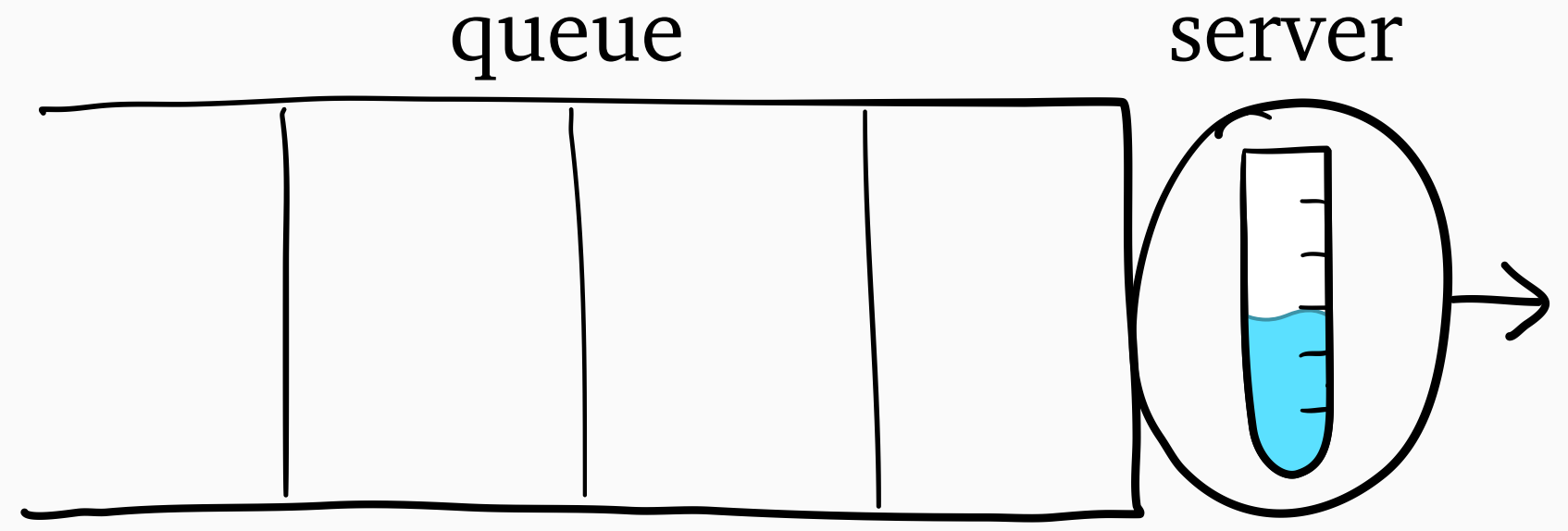
# Single-server queueing system

M/G/1



# Single-server queueing system

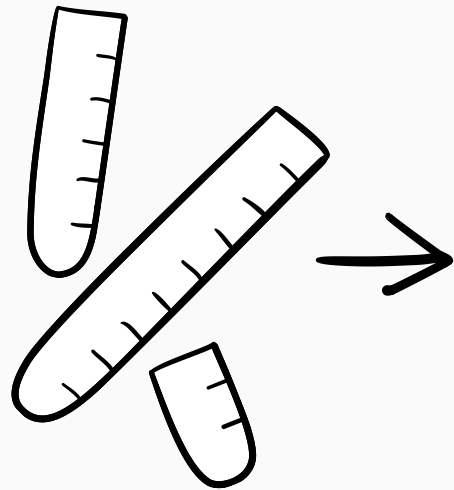
M/G/1



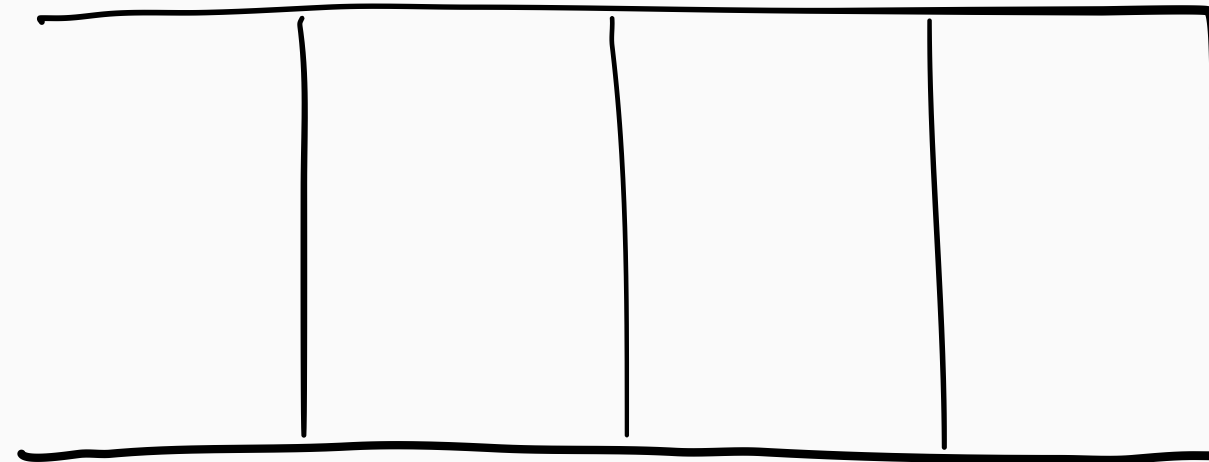
# Single-server queueing system

M/G/1

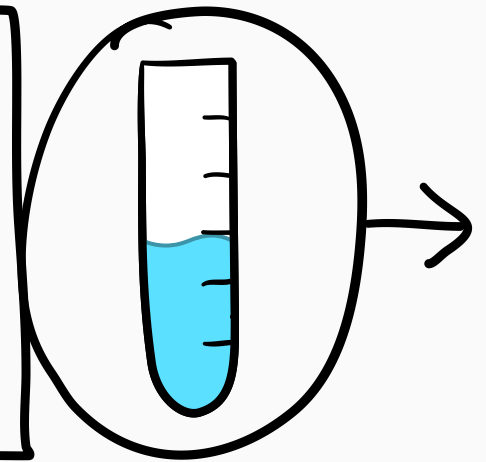
stochastic arrivals



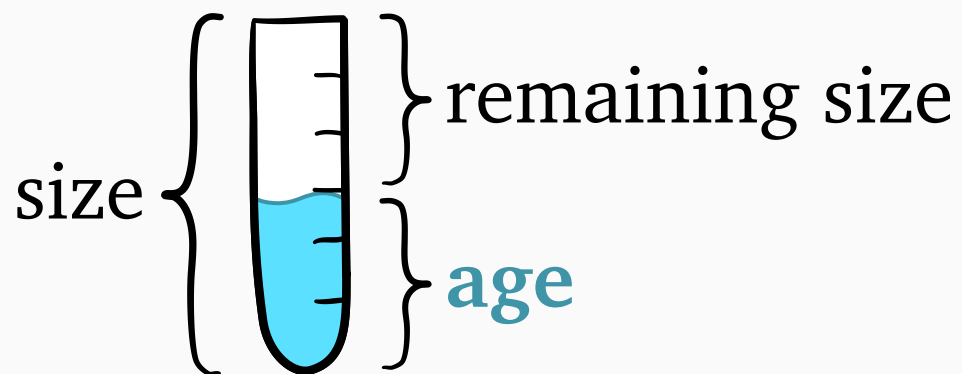
queue



server



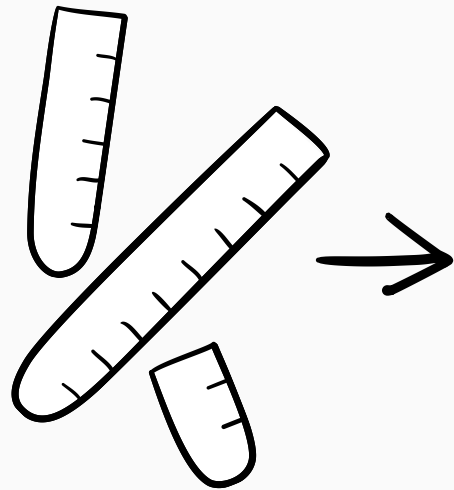
job



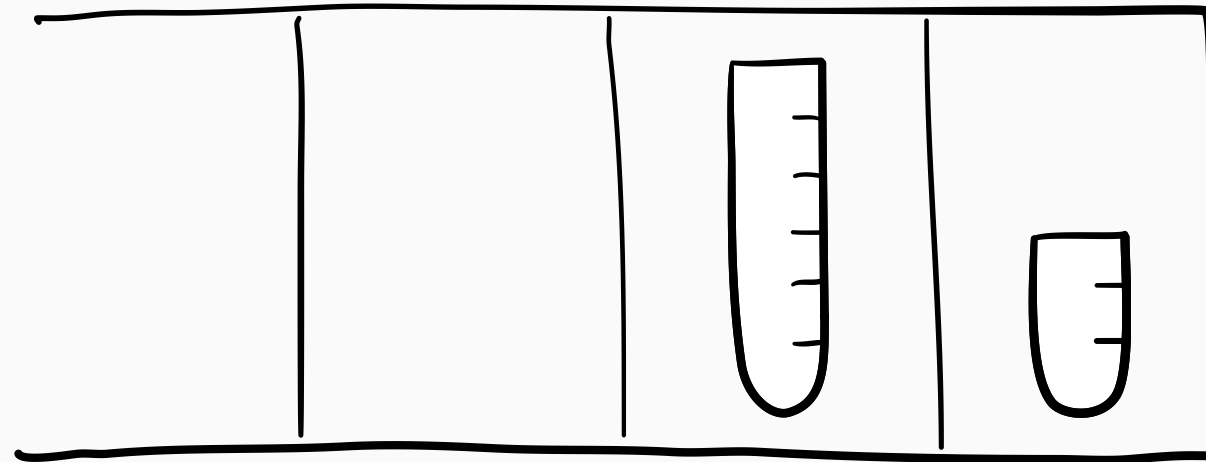
# Single-server queueing system

M/G/1

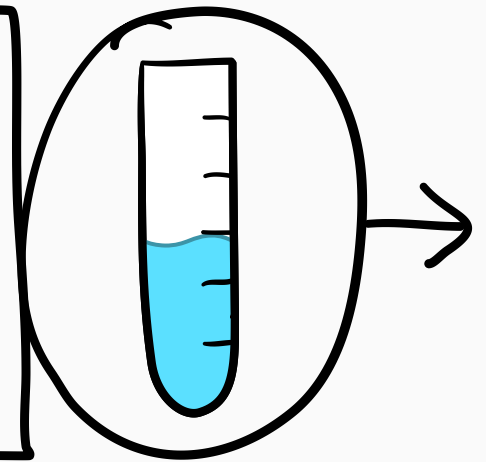
stochastic arrivals



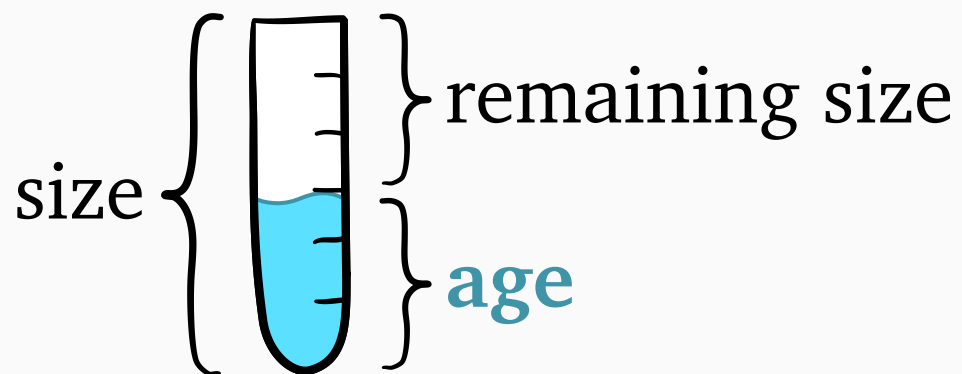
queue



server



job



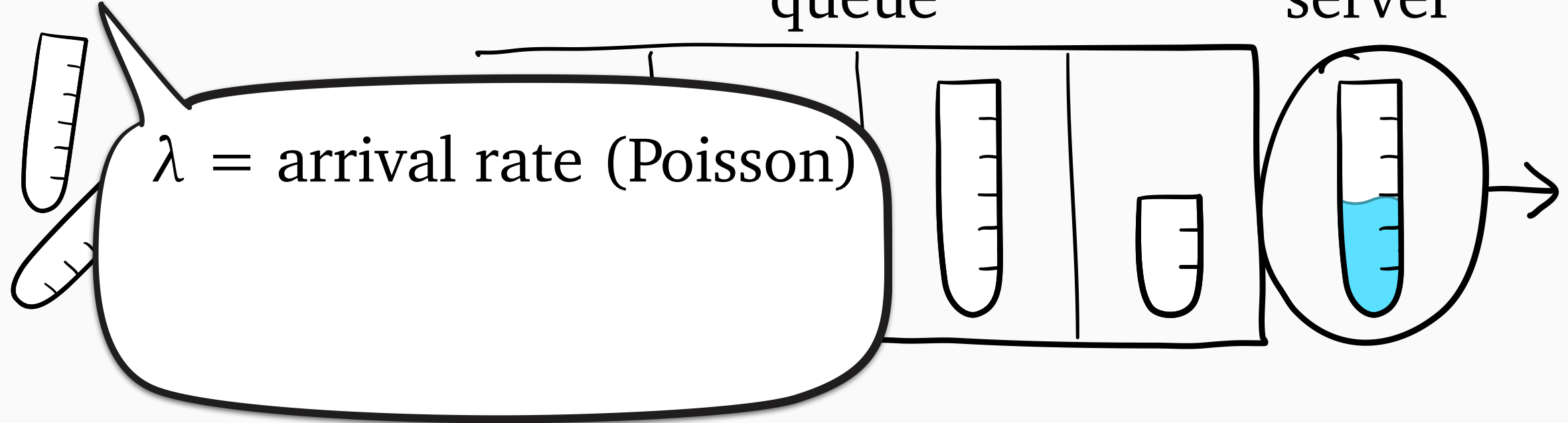
# Single-server queueing system

M/G/1

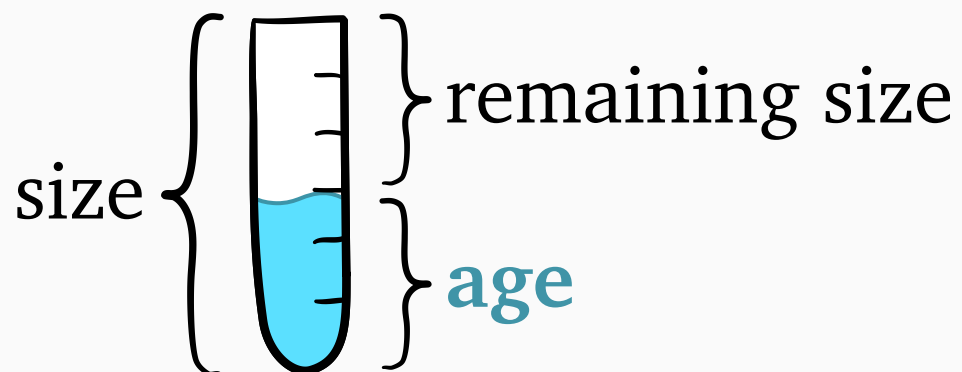
stochastic arrivals

queue

server



job



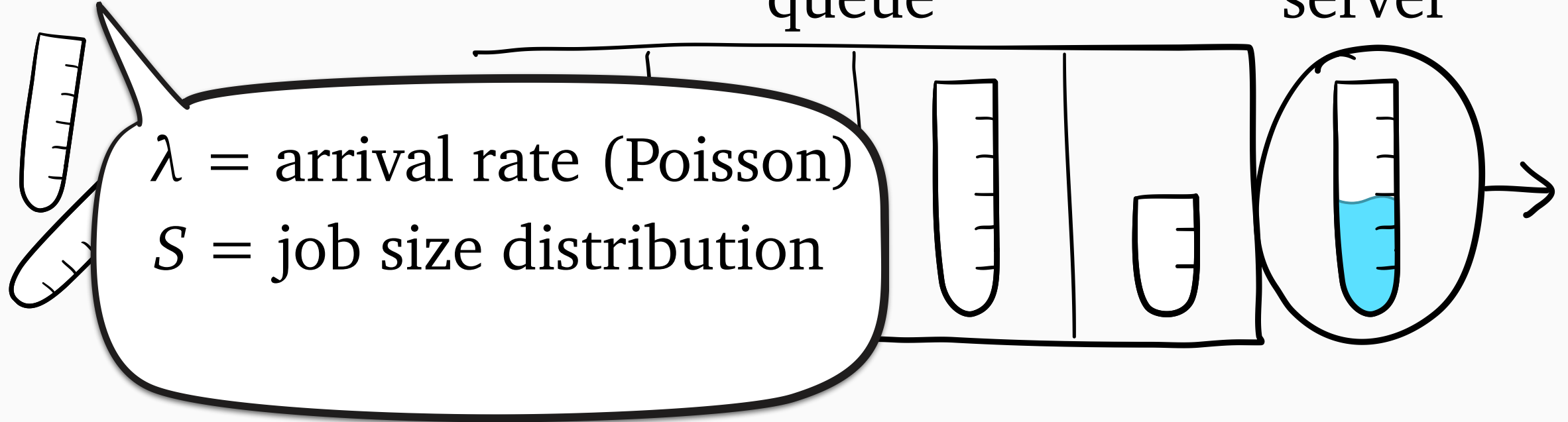
# Single-server queueing system

M/G/1

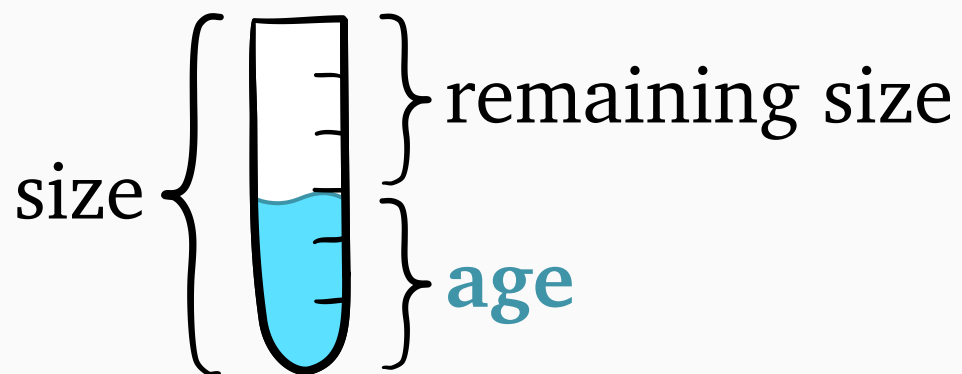
stochastic arrivals

queue

server



job



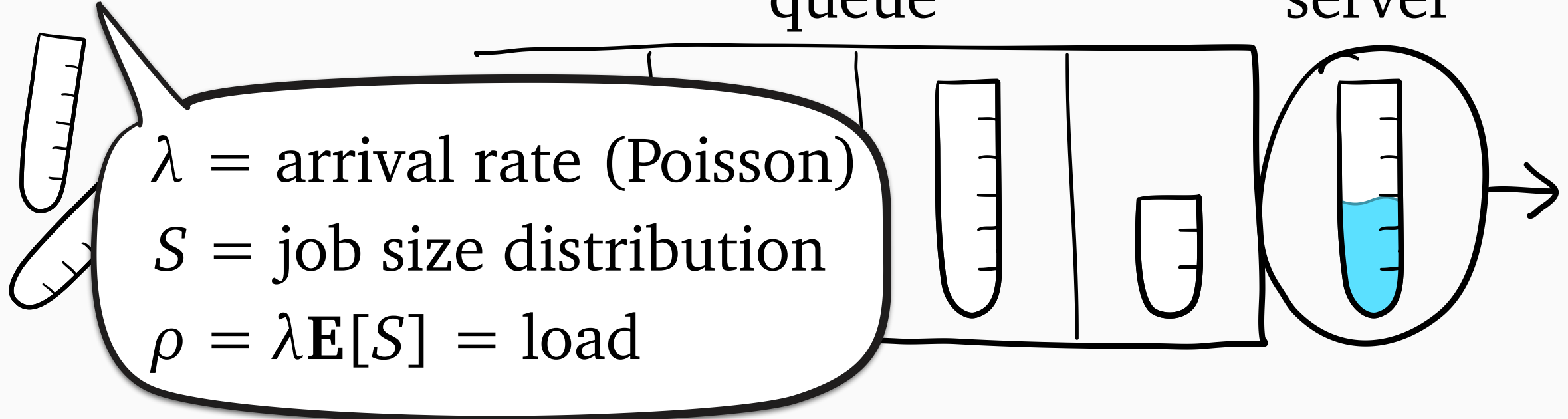
# Single-server queueing system

M/G/1

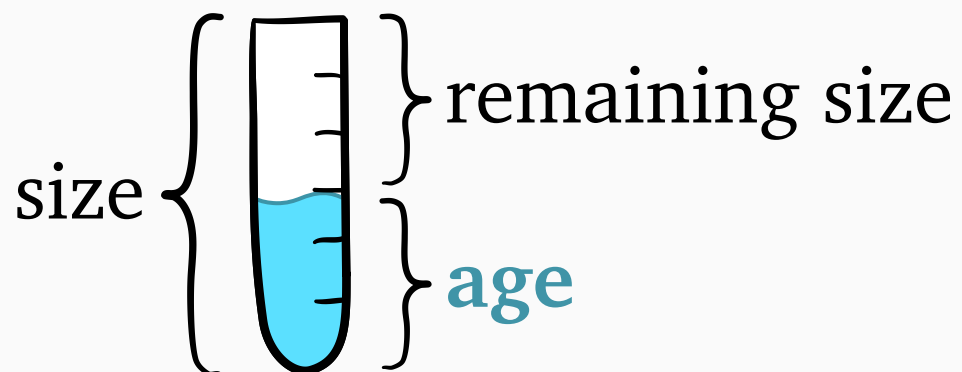
stochastic arrivals

queue

server



job





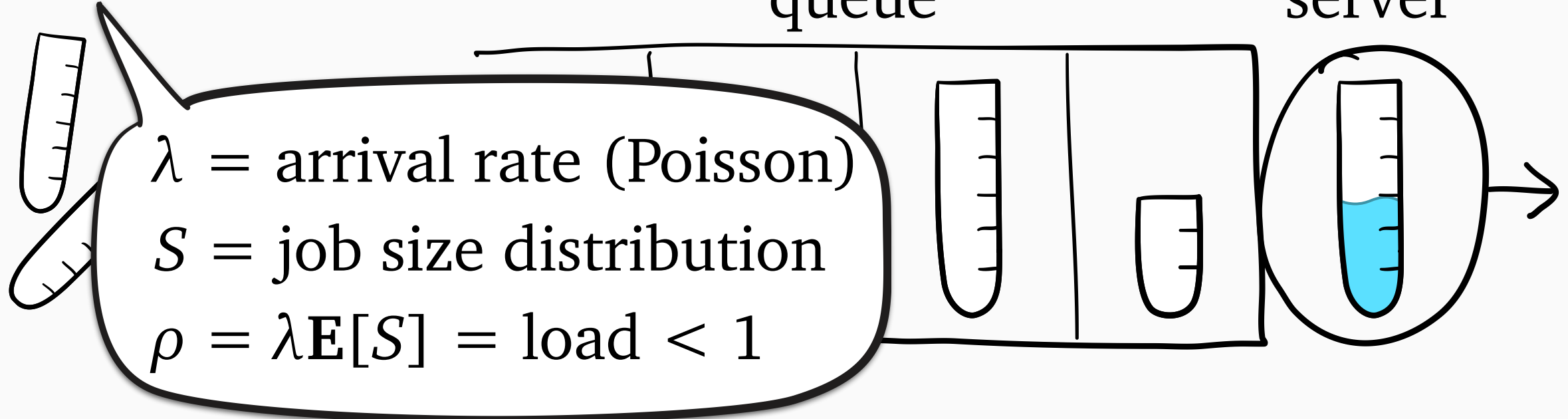
# Single-server queueing system

M/G/1

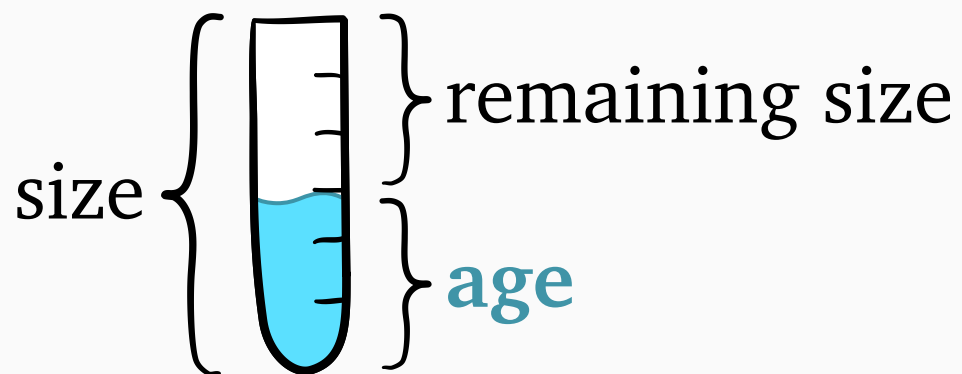
stochastic arrivals

queue

server



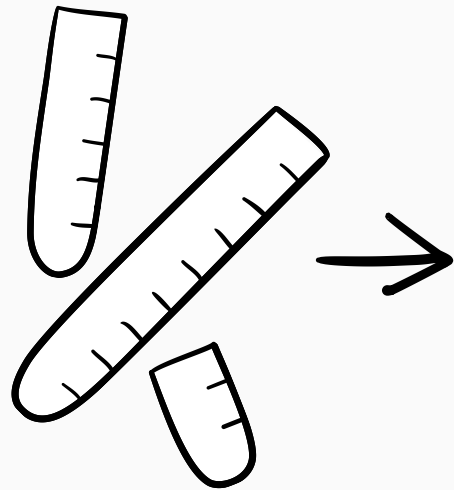
job



# Single-server queueing system

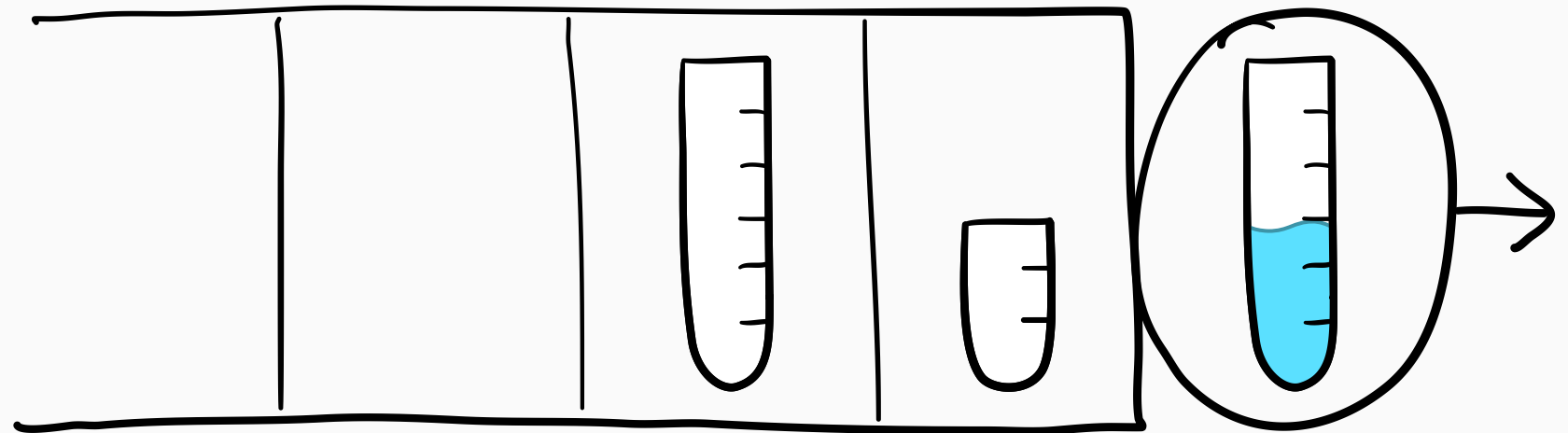
M/G/1

stochastic arrivals

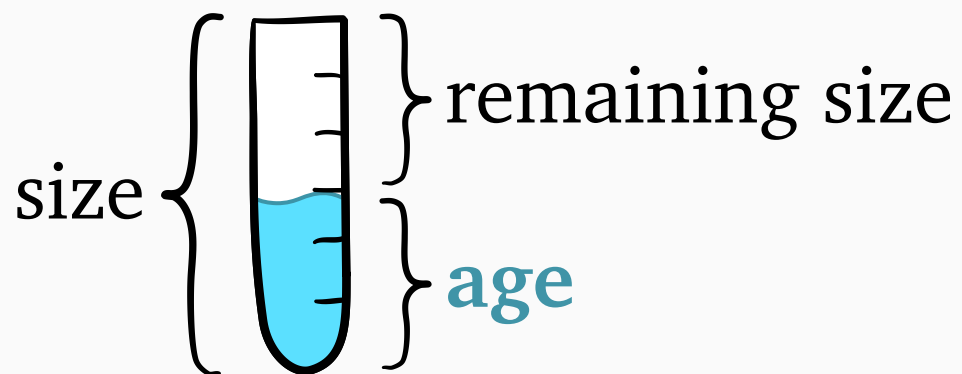


queue

server



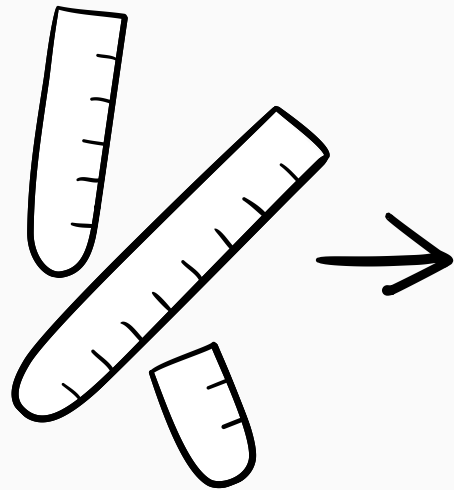
job



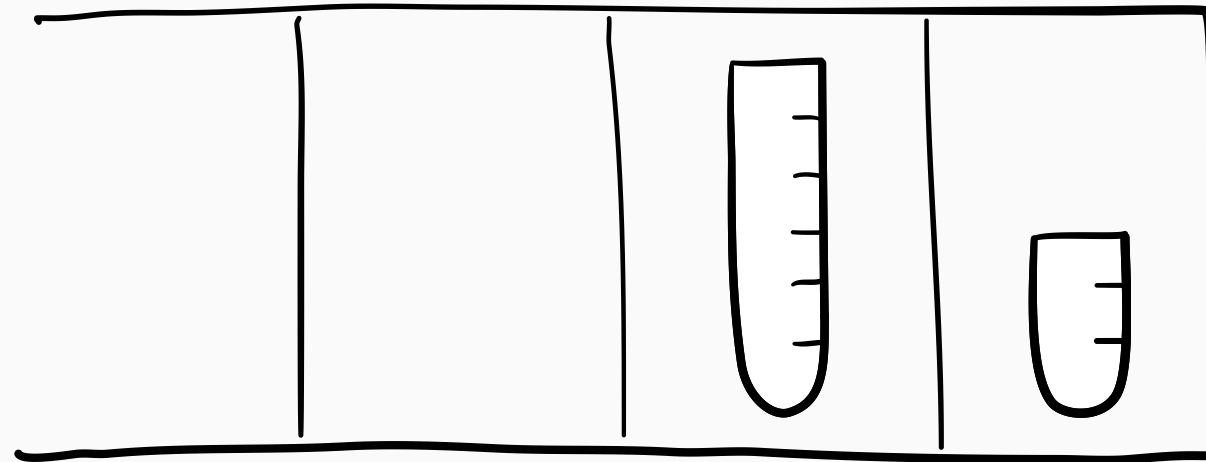
# Single-server queueing system

M/G/1

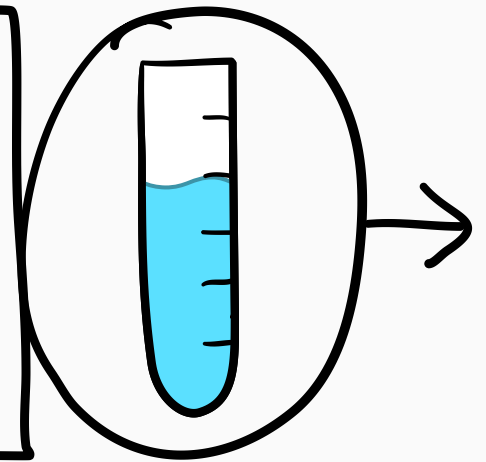
stochastic arrivals



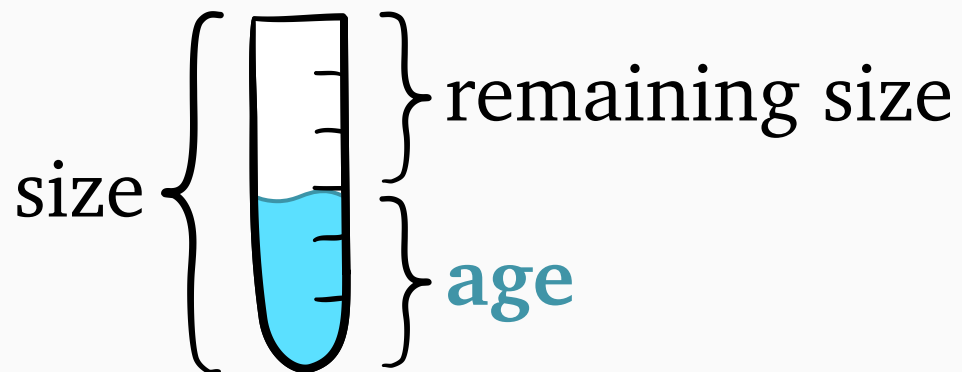
queue



server



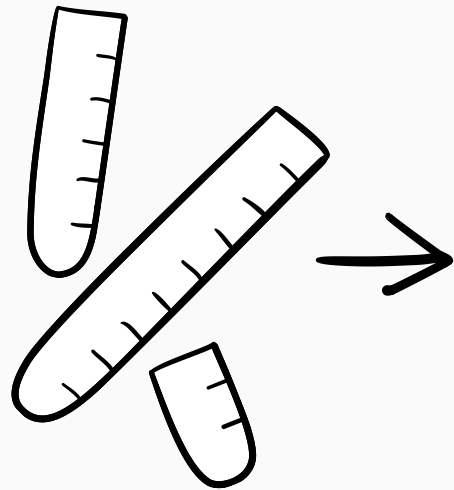
job



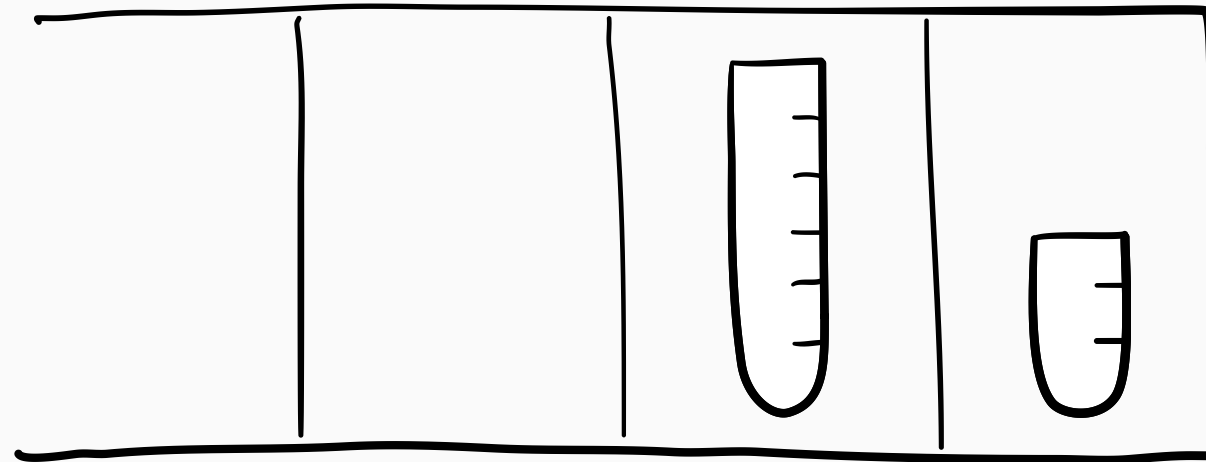
# Single-server queueing system

M/G/1

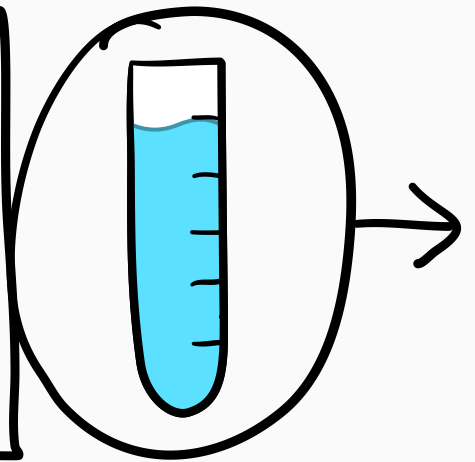
stochastic arrivals



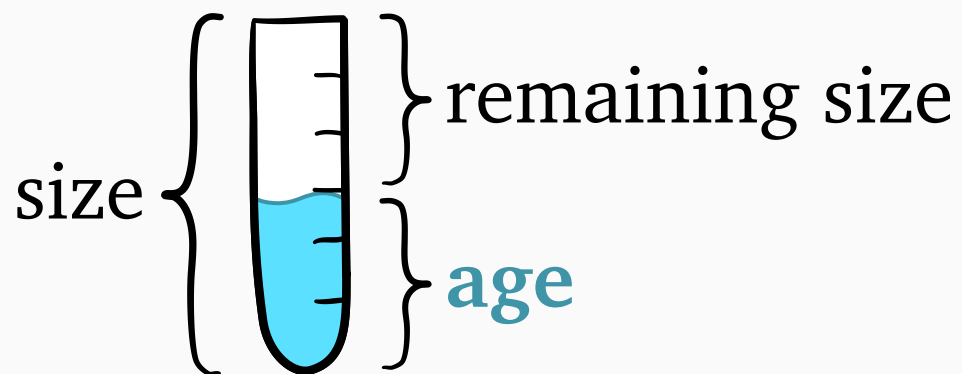
queue



server



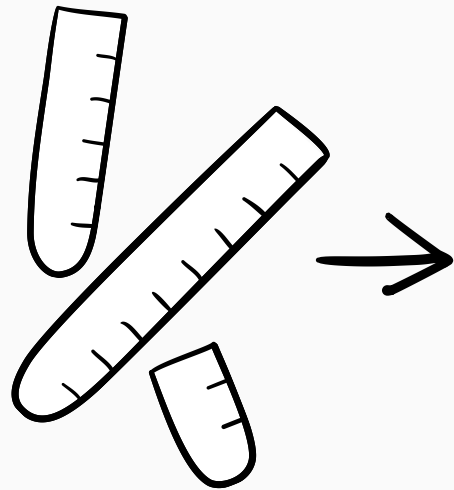
job



# Single-server queueing system

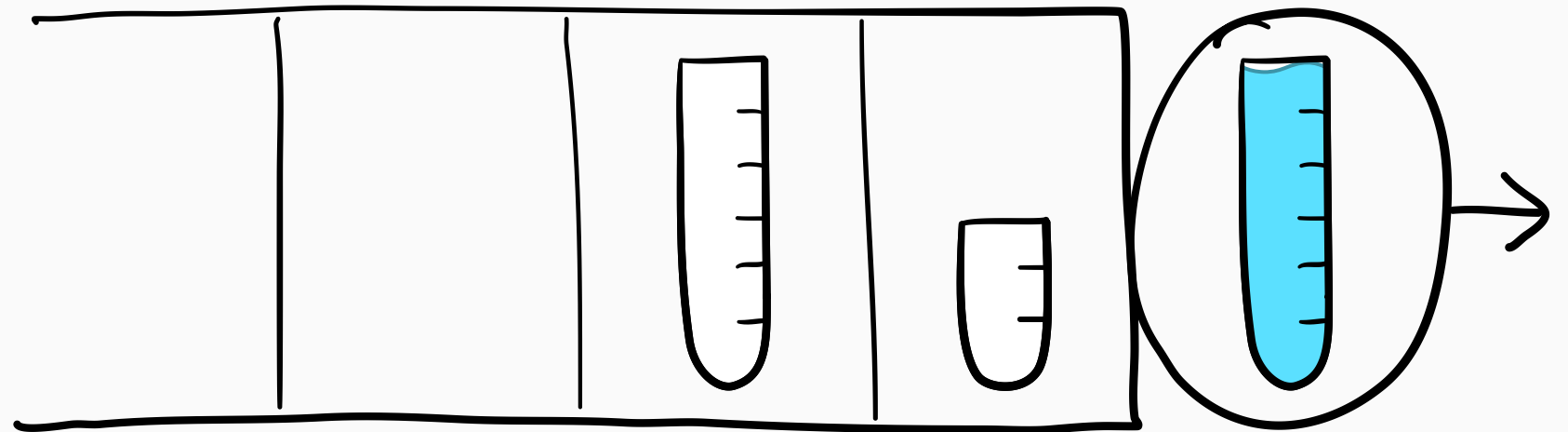
M/G/1

stochastic arrivals

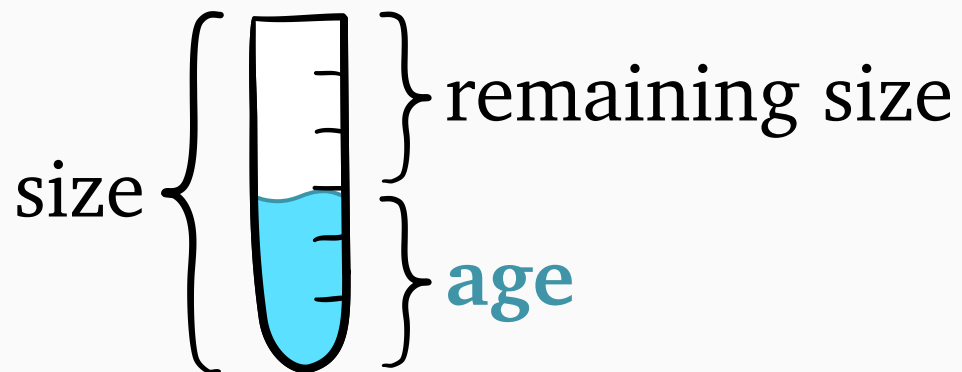


queue

server



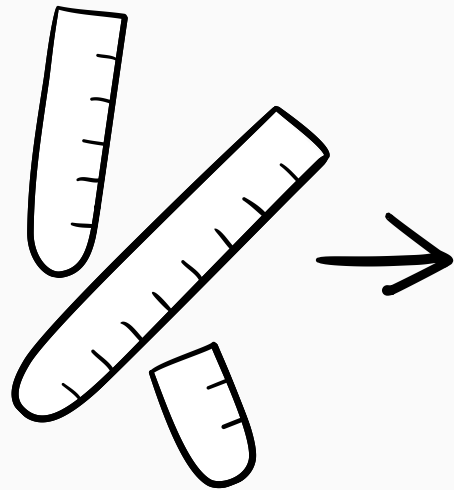
job



# Single-server queueing system

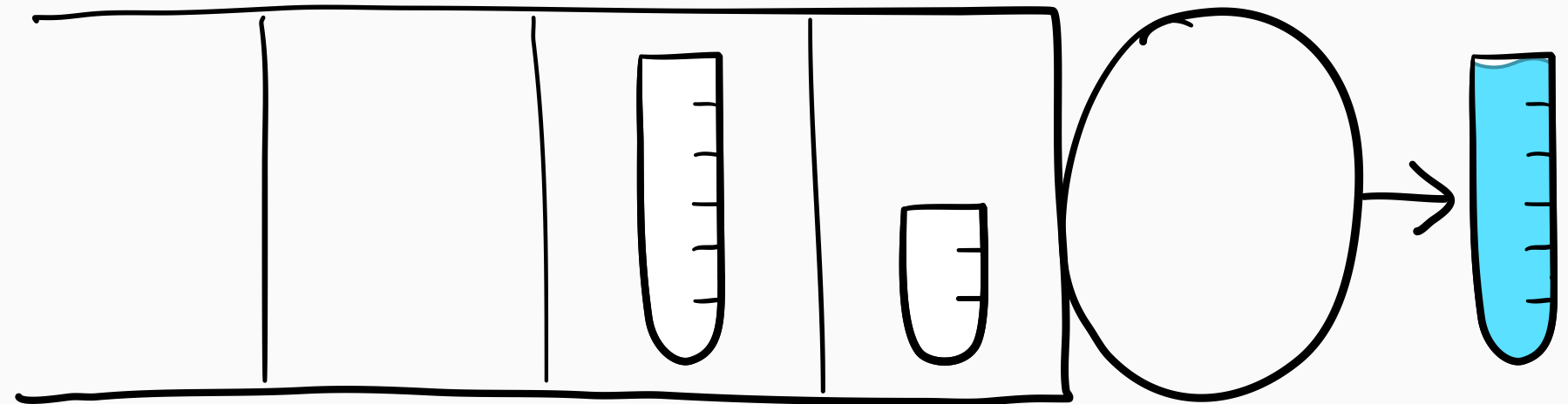
M/G/1

stochastic arrivals

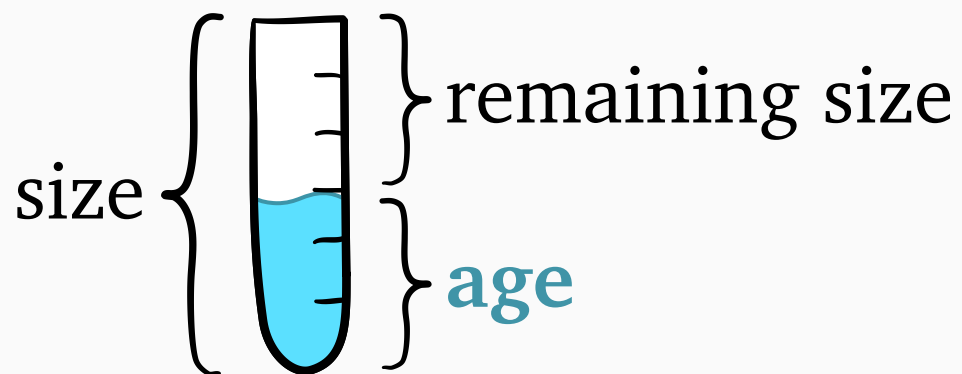


queue

server



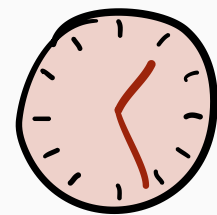
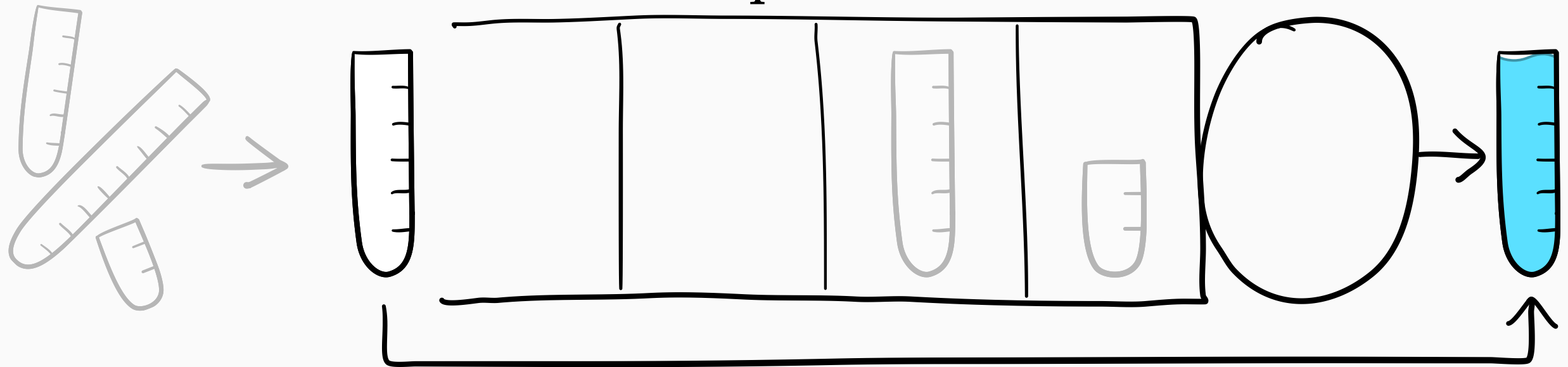
job



# Single-server queueing system

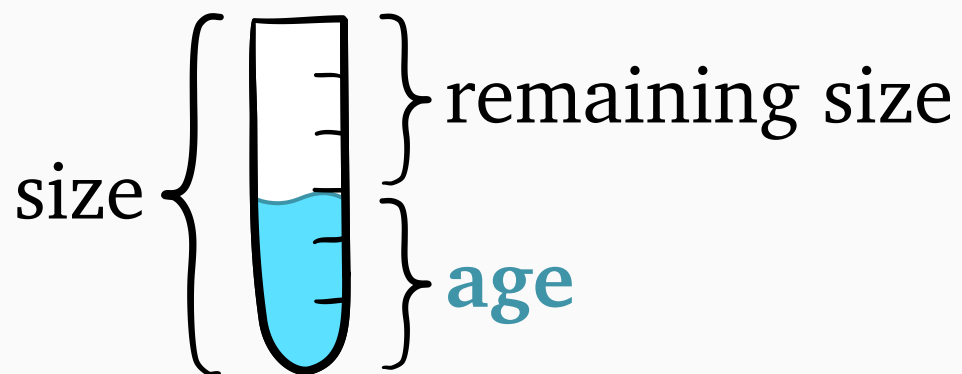
M/G/1

stochastic arrivals



response time

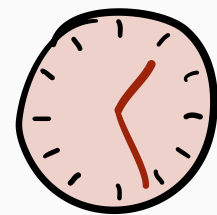
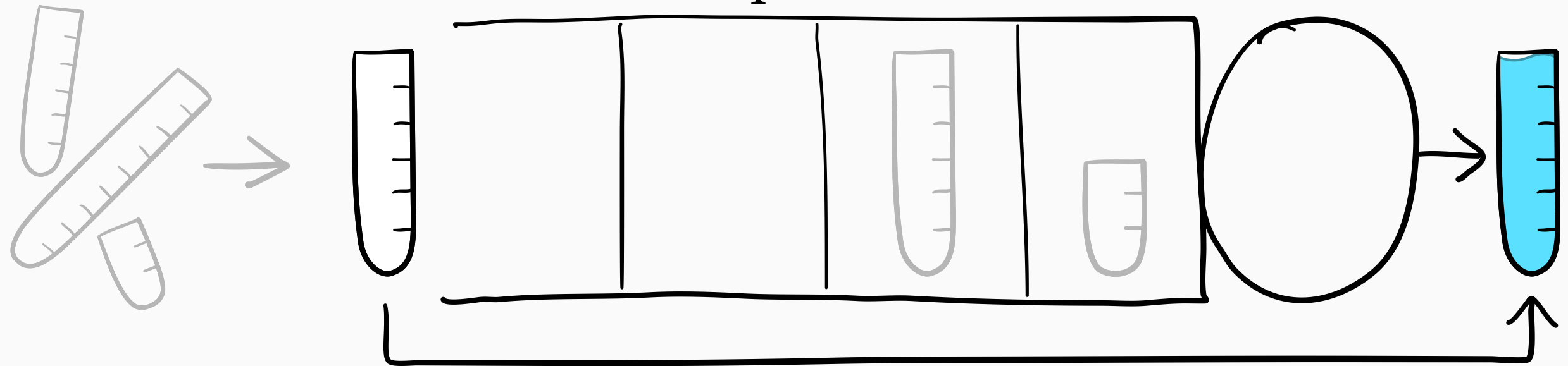
job



# Single-server queueing system

M/G/1

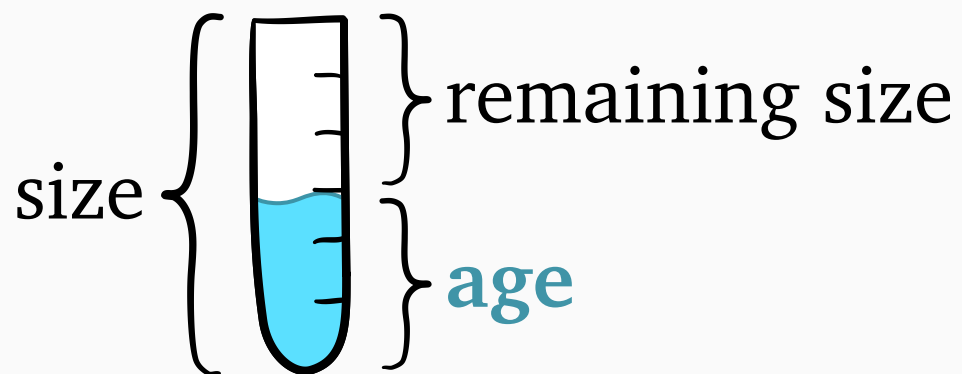
stochastic arrivals



response time

a.k.a. delay, latency

job

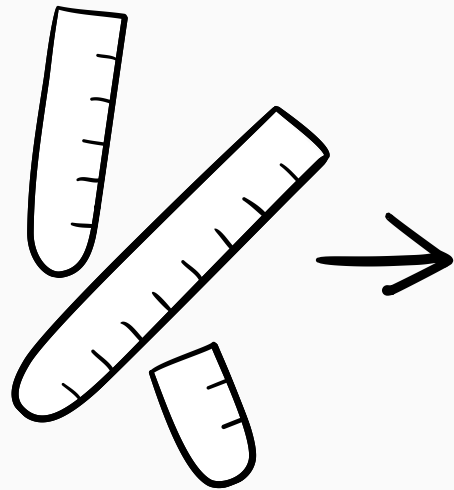




# Single-server queueing system

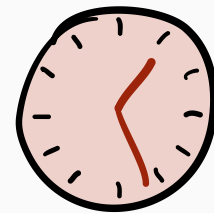
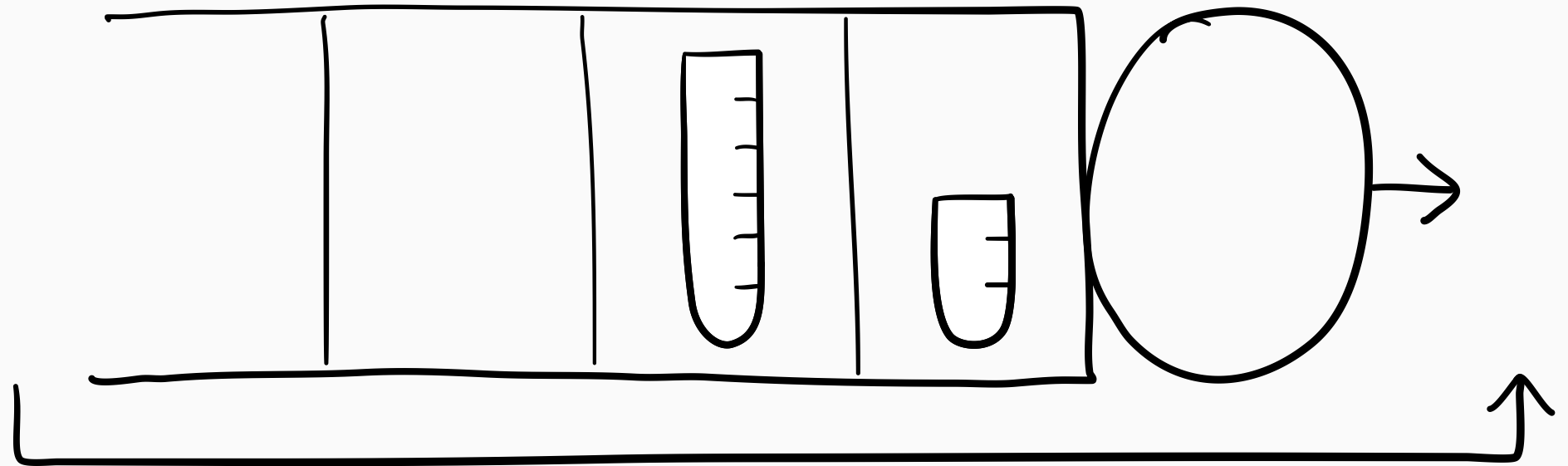
M/G/1

stochastic arrivals



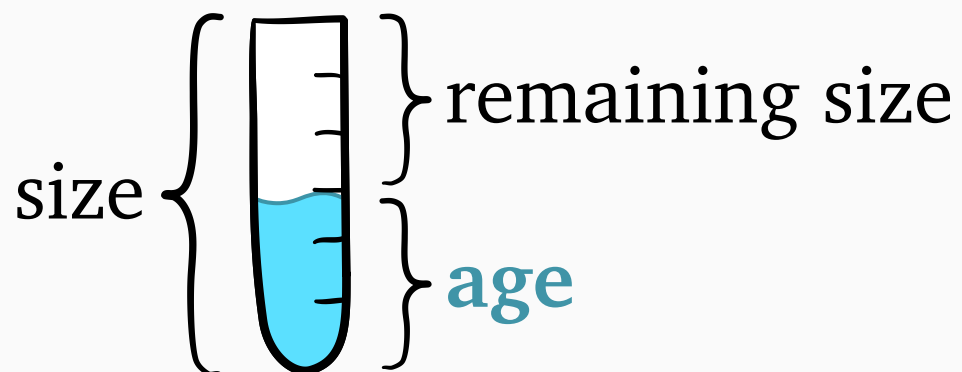
queue

server



response time

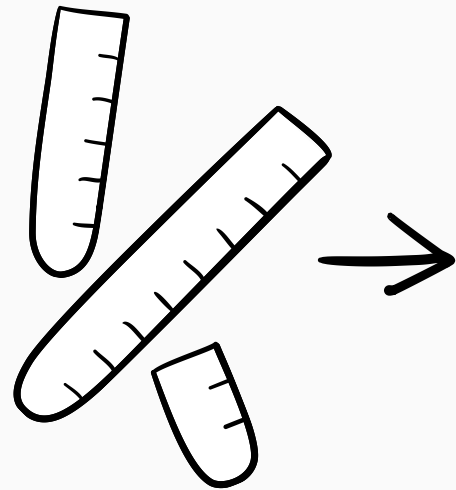
job



# Single-server queueing system

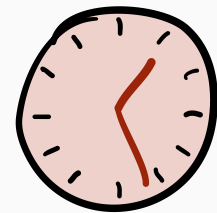
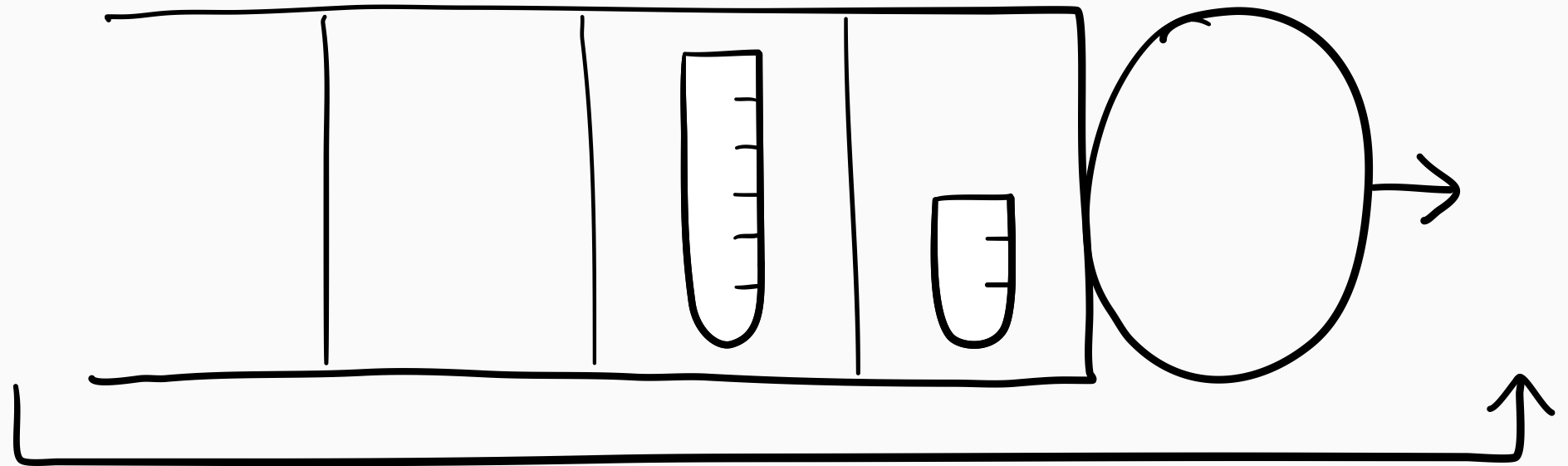
M/G/1

stochastic arrivals



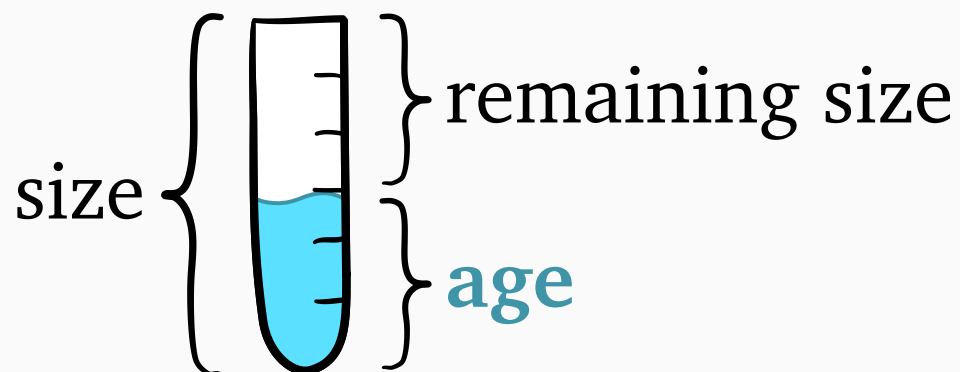
queue

server



response time

job

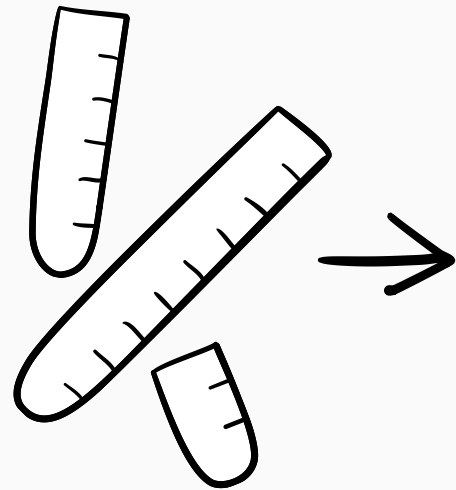


**Scheduling policy:**  
decides which job to serve

# Single-server queueing system

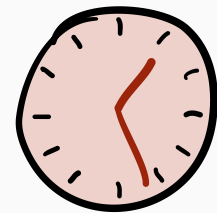
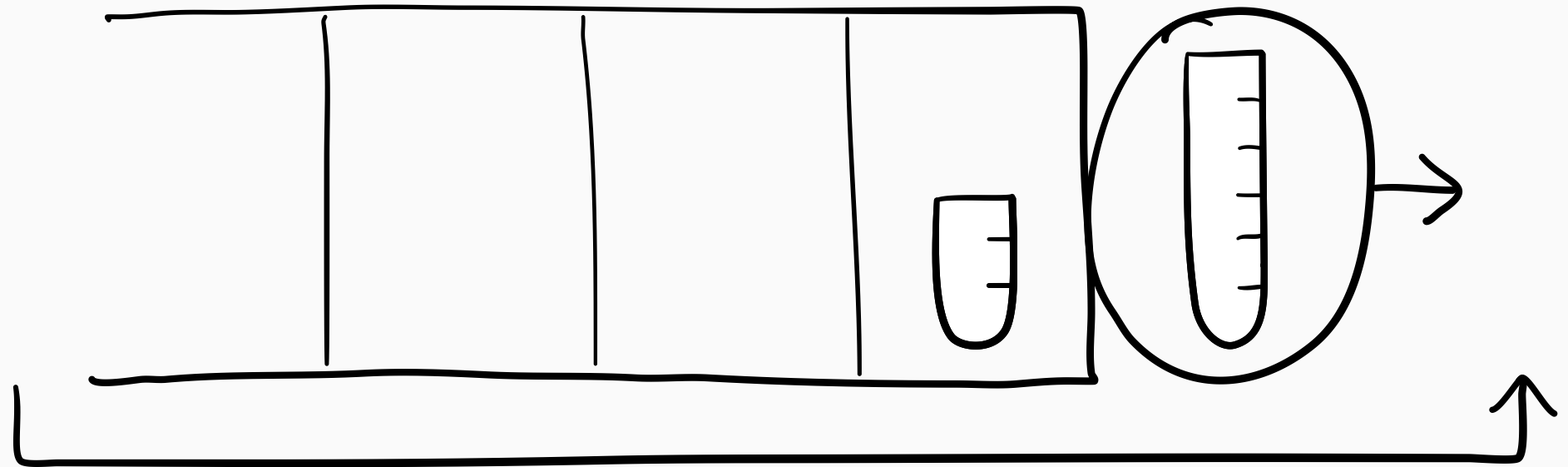
M/G/1

stochastic arrivals



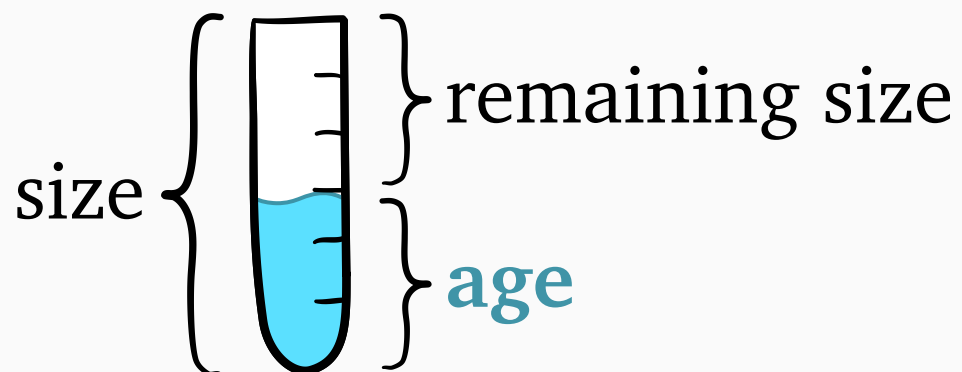
queue

server



response time

job

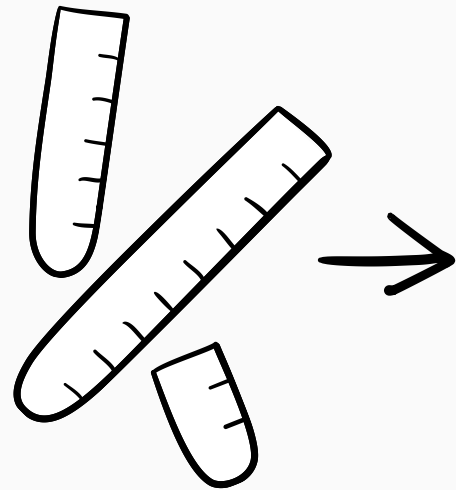


**Scheduling policy:**  
decides which job to serve

# Single-server queueing system

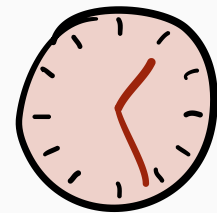
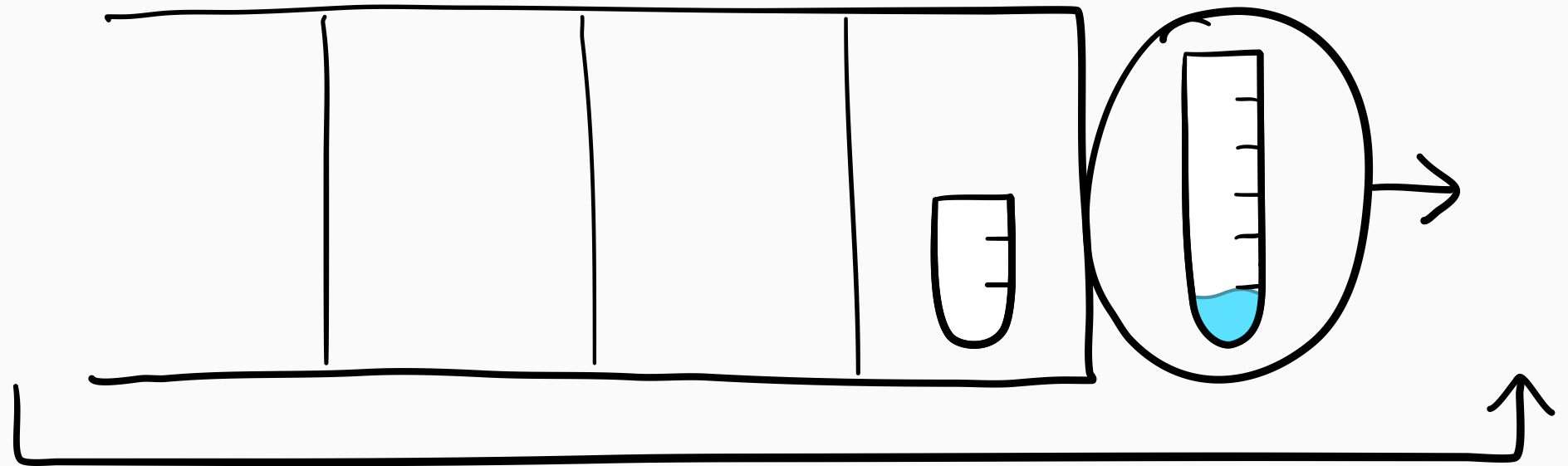
M/G/1

stochastic arrivals



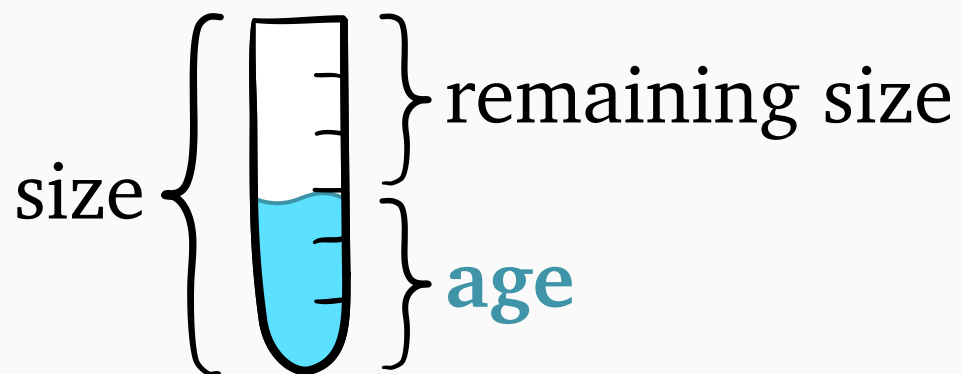
queue

server



response time

job

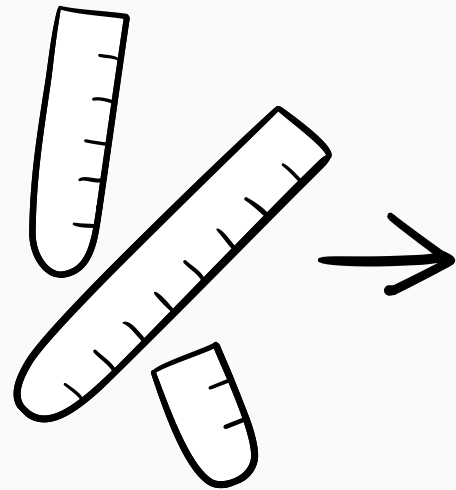


**Scheduling policy:**  
decides which job to serve

# Single-server queueing system

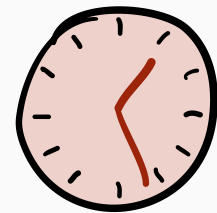
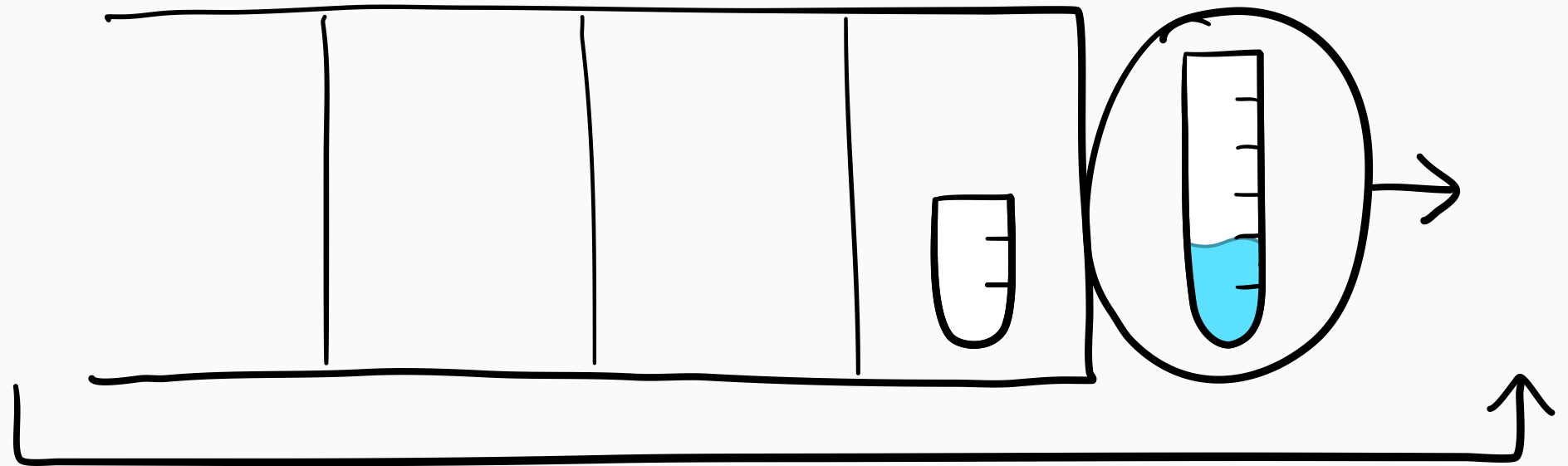
M/G/1

stochastic arrivals



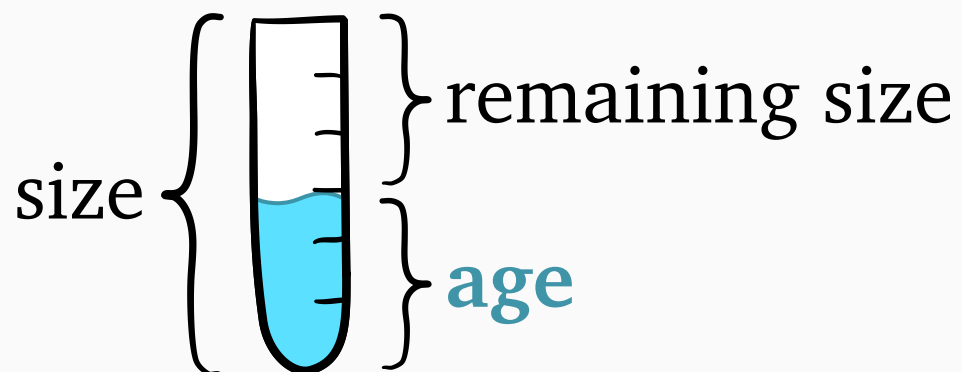
queue

server



response time

job

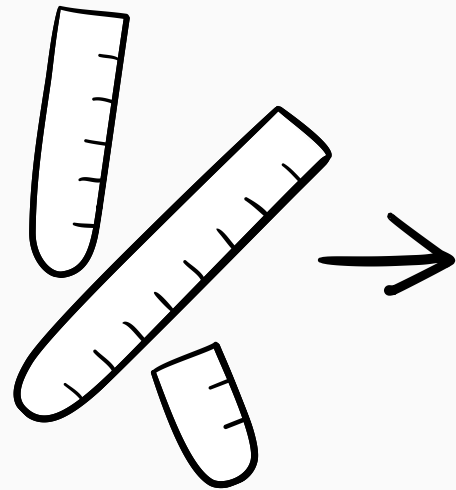


**Scheduling policy:**  
decides which job to serve

# Single-server queueing system

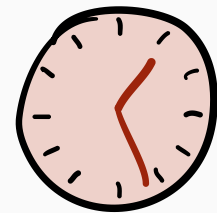
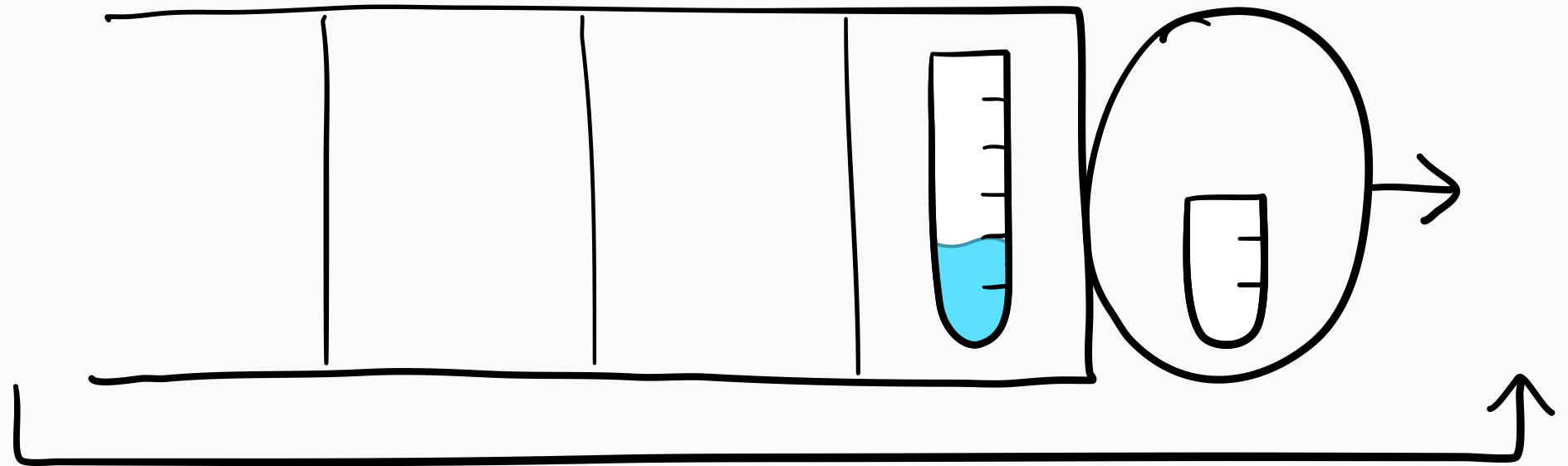
M/G/1

stochastic arrivals



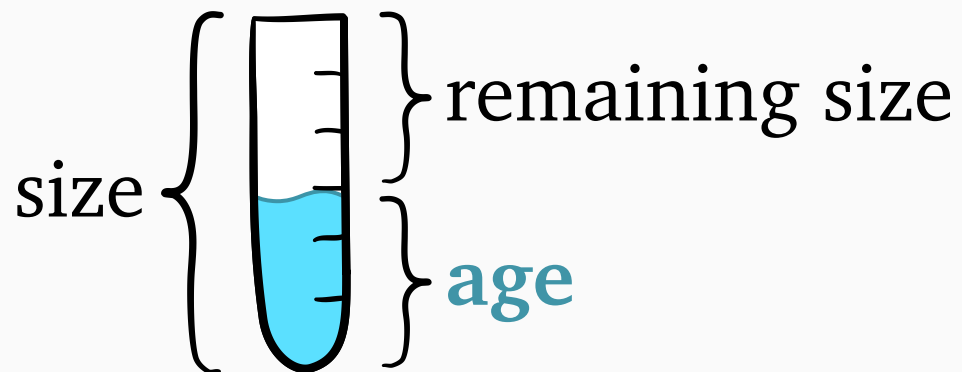
queue

server



response time

job

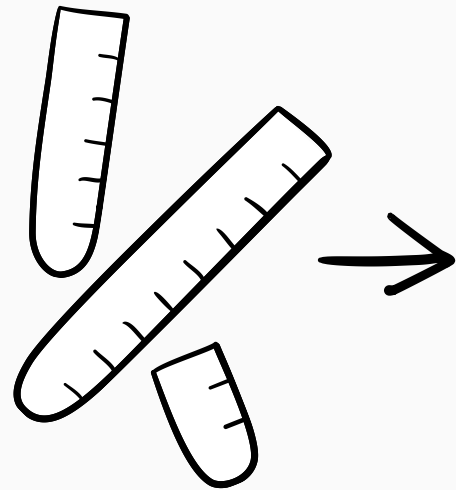


**Scheduling policy:**  
decides which job to serve

# Single-server queueing system

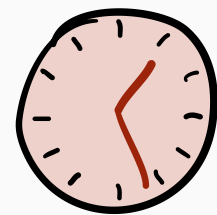
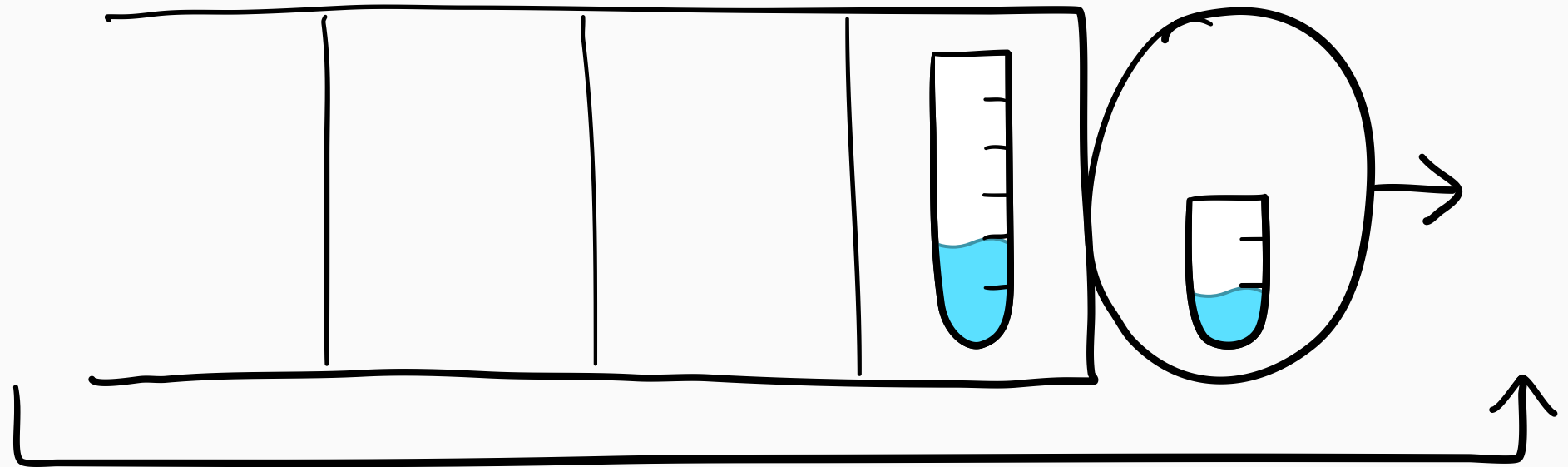
M/G/1

stochastic arrivals



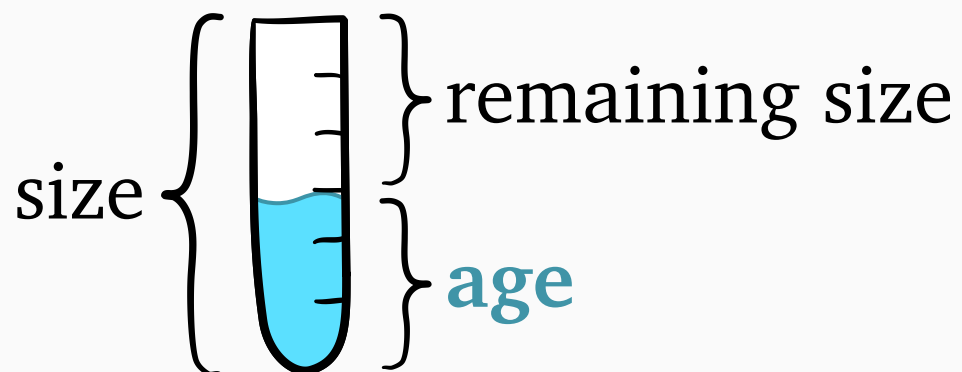
queue

server



response time

job

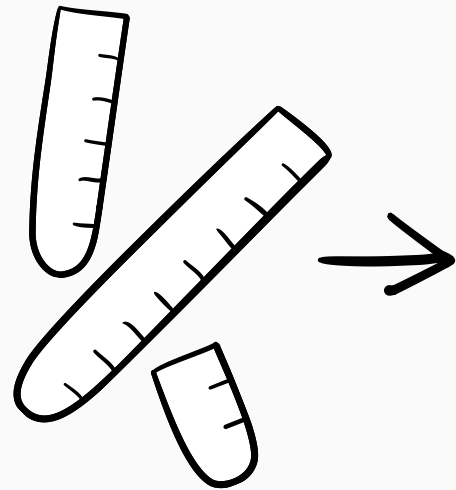


**Scheduling policy:**  
decides which job to serve

# Single-server queueing system

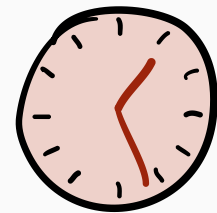
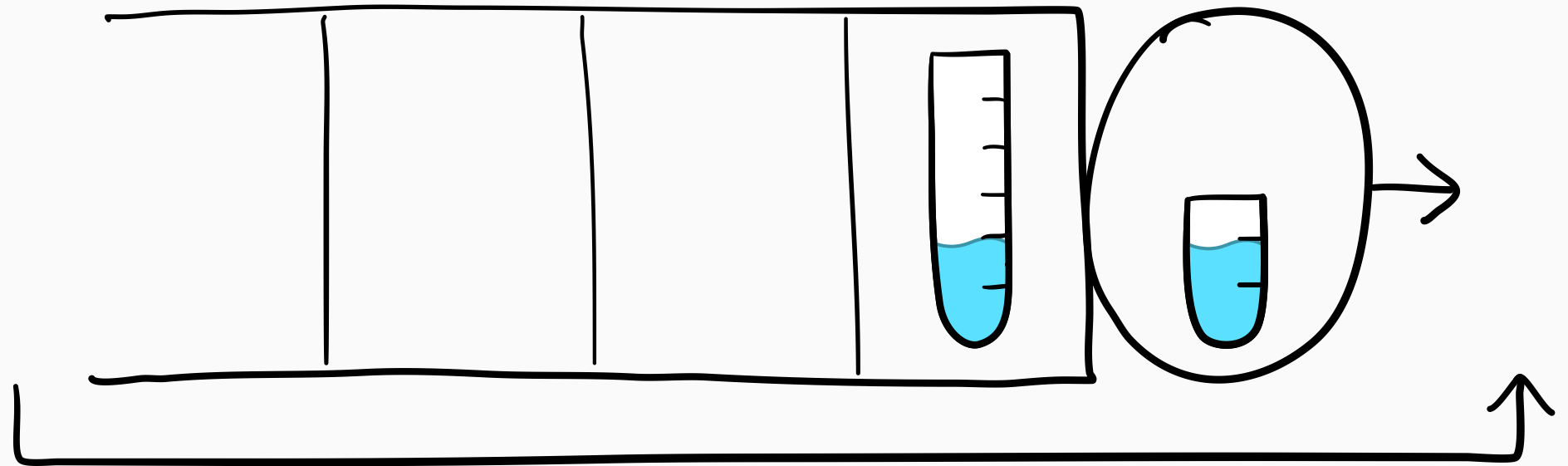
M/G/1

stochastic arrivals



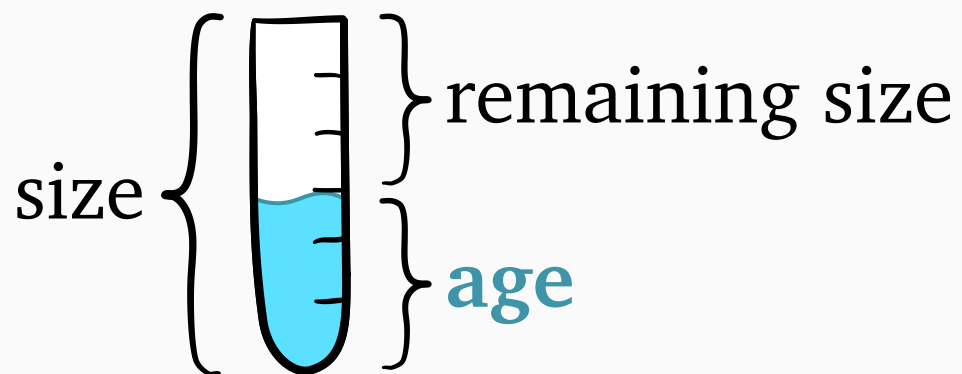
queue

server



response time

job



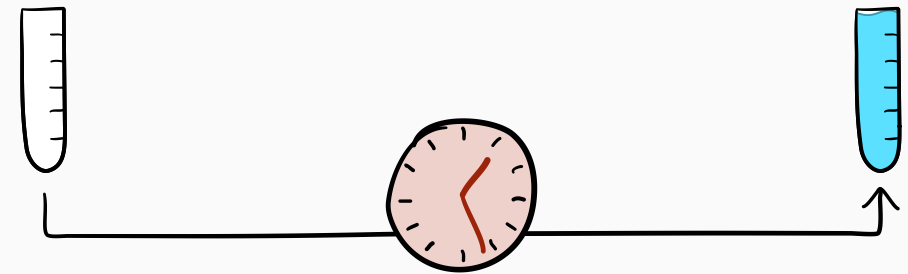
**Scheduling policy:**  
decides which job to serve



# Evaluating performance

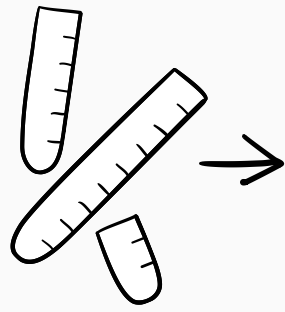
# Evaluating performance

response time

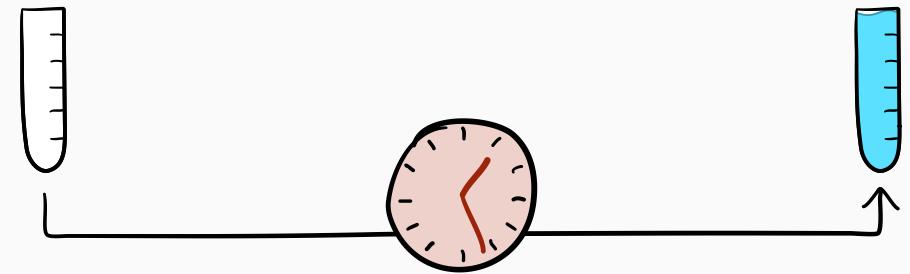


# Evaluating performance

stochastic arrival  
process  $\lambda, S$

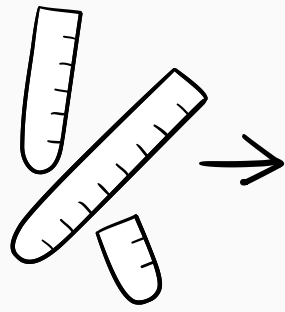


response time

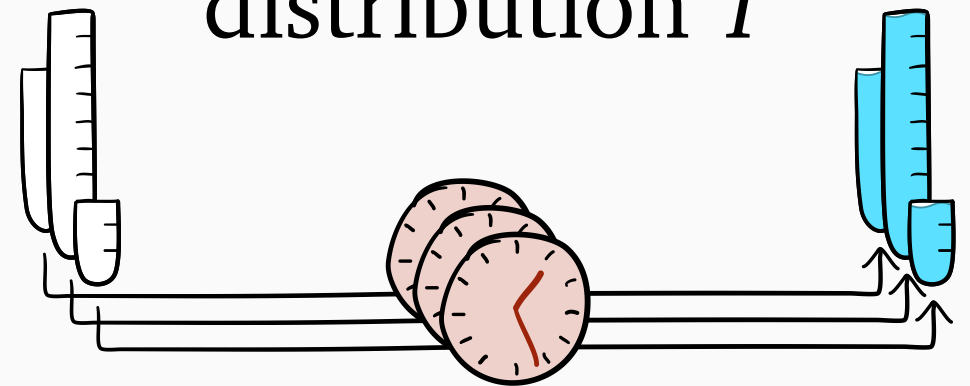


# Evaluating performance

stochastic arrival  
process  $\lambda, S$

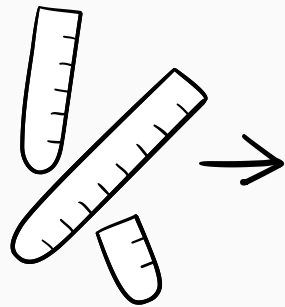


response time  
distribution  $T$



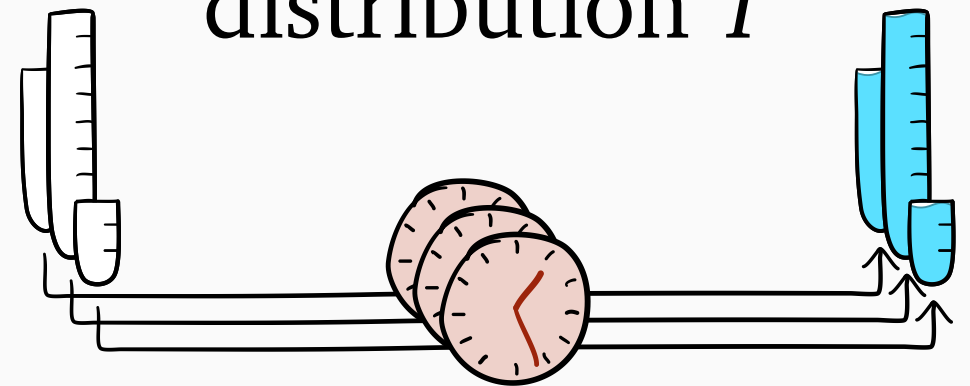
# Evaluating performance

stochastic arrival  
process  $\lambda, S$



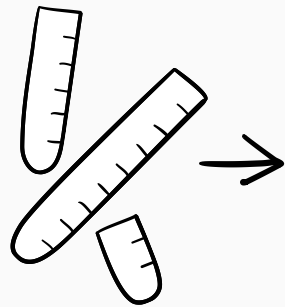
**scheduling  
policy**

response time  
distribution  $T$



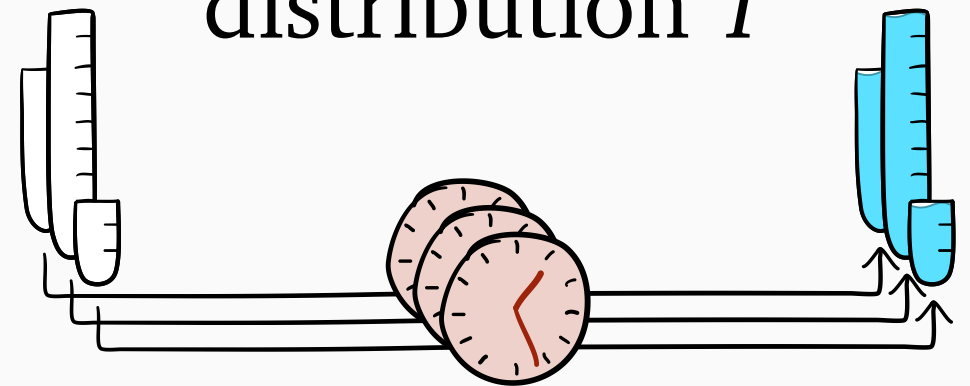
# Evaluating performance

stochastic arrival  
process  $\lambda, S$



scheduling  
policy

response time  
distribution  $T$

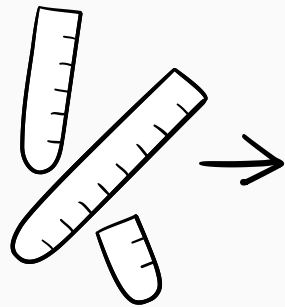


## Metrics:

- mean  $E[T]$
- tail probability  $P[T > t]$
- percentiles, e.g.  $T_{99}$

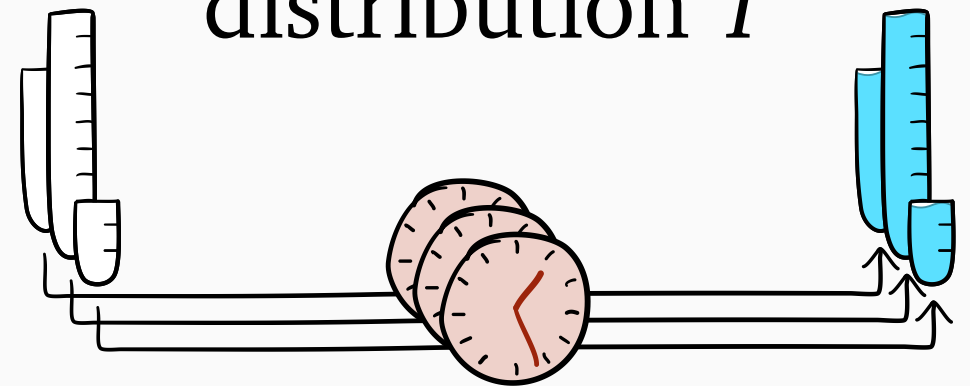
# Evaluating performance

stochastic arrival  
process  $\lambda, S$



scheduling  
policy

response time  
distribution  $T$



**Goals:**

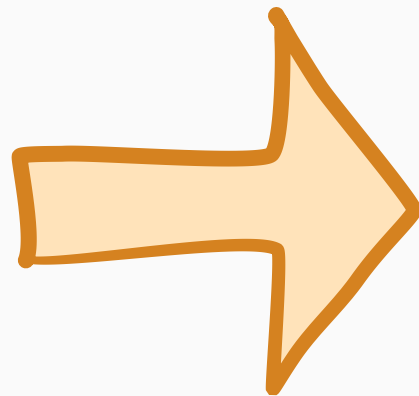
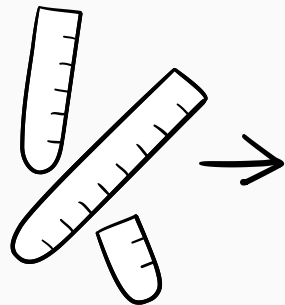
- analyze **policy**

**Metrics:**

- mean  $E[T]$
- tail probability  $P[T > t]$
- percentiles, e.g.  $T_{99}$

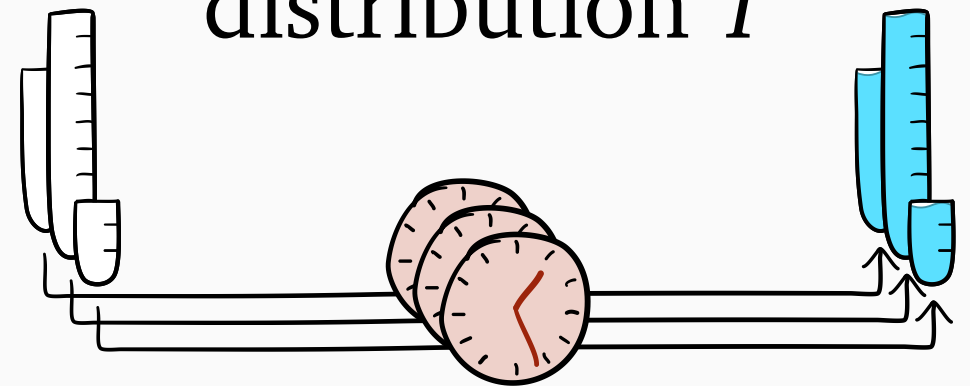
# Evaluating performance

stochastic arrival  
process  $\lambda, S$



**scheduling  
policy**

response time  
distribution  $T$



**Goals:**

- *analyze* **policy**
- *optimize* **policy**

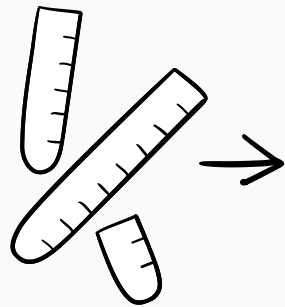
**Metrics:**

- mean  $E[T]$
- tail probability  $\mathbf{P}[T > t]$
- percentiles, e.g.  $T_{99}$



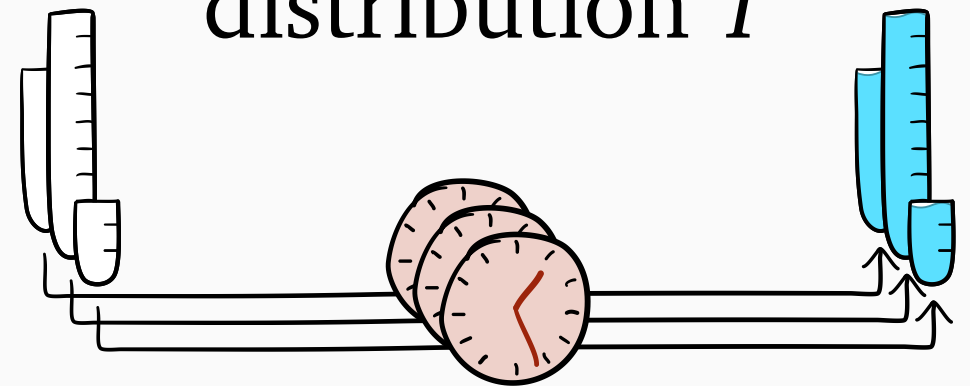
# Evaluating performance

stochastic arrival  
process  $\lambda, S$



scheduling  
policy

response time  
distribution  $T$



**Goals:**

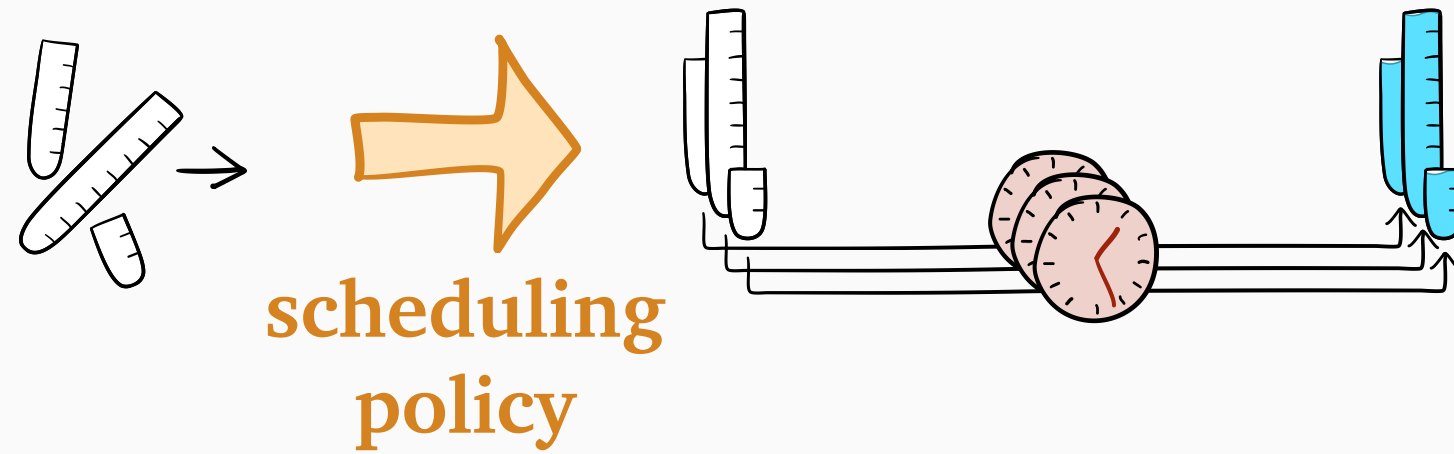
our focus

- analyze **policy**
- optimize **policy**

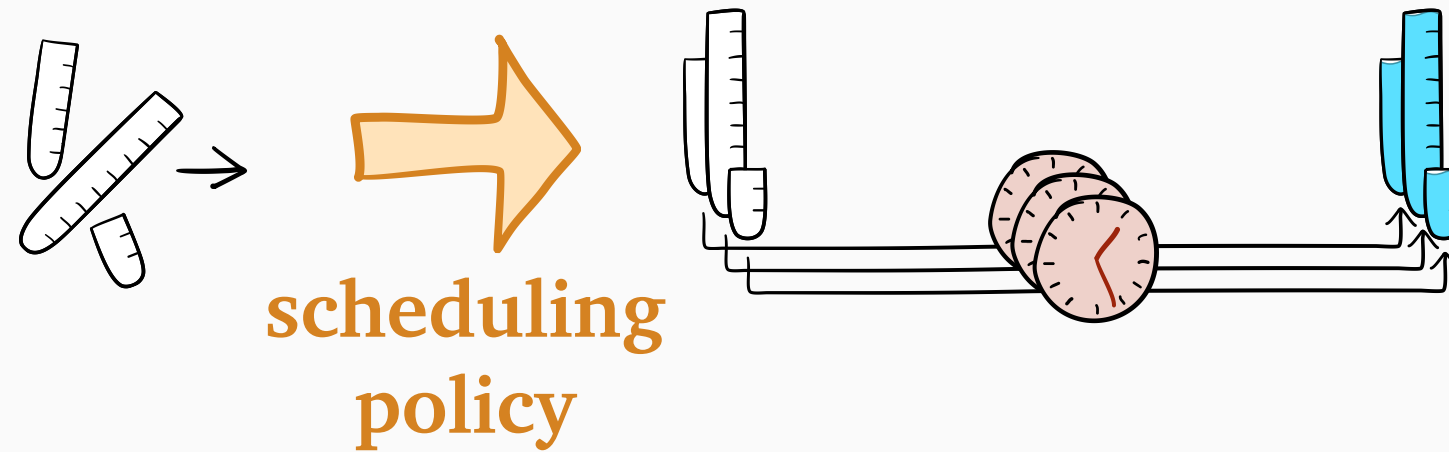
**Metrics:**

- mean  $E[T]$
- tail probability  $P[T > t]$
- percentiles, e.g.  $T_{99}$

# Goal: analyze impact of **policy**



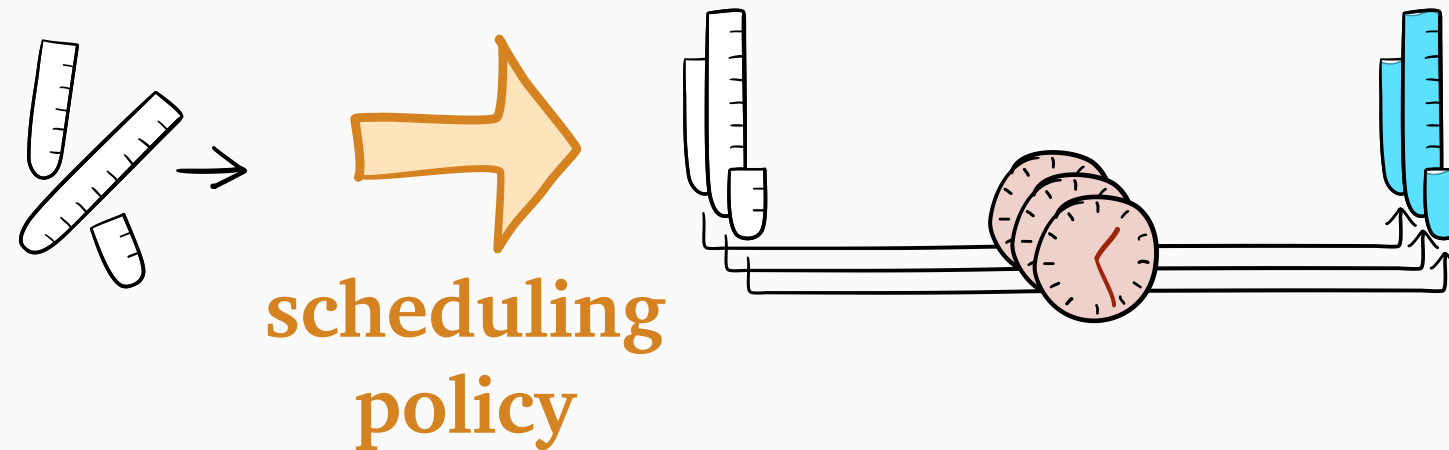
# Goal: analyze impact of **policy**



## Problem:

- Theory limited to “simple” **policies**

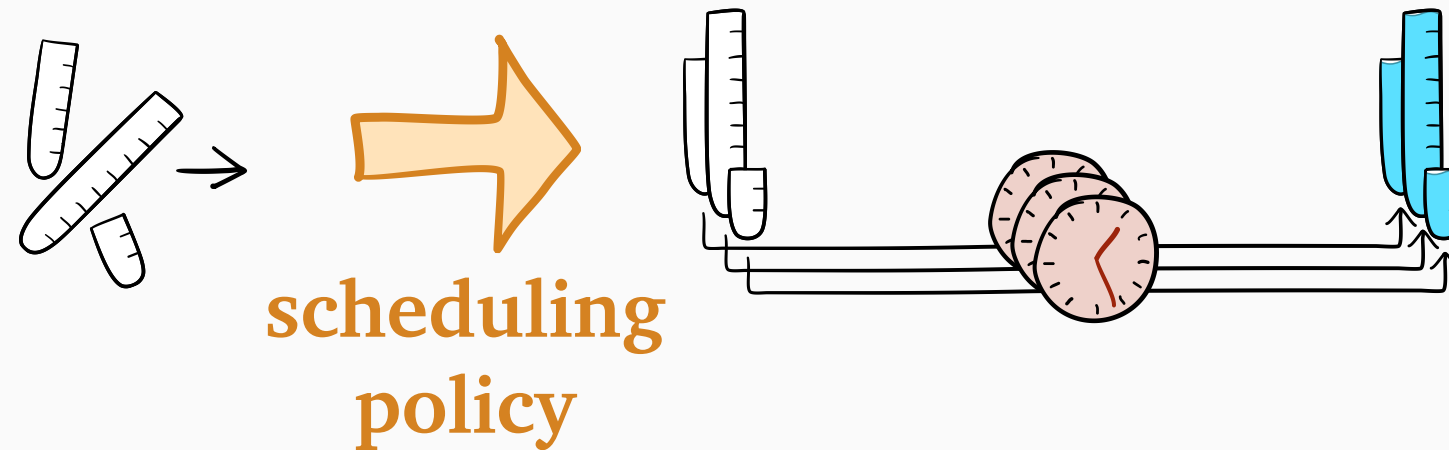
# Goal: analyze impact of **policy**



## Problem:

- Theory limited to “simple” **policies**
- In practice, “complex” **policies** matter

# Goal: analyze impact of **policy**



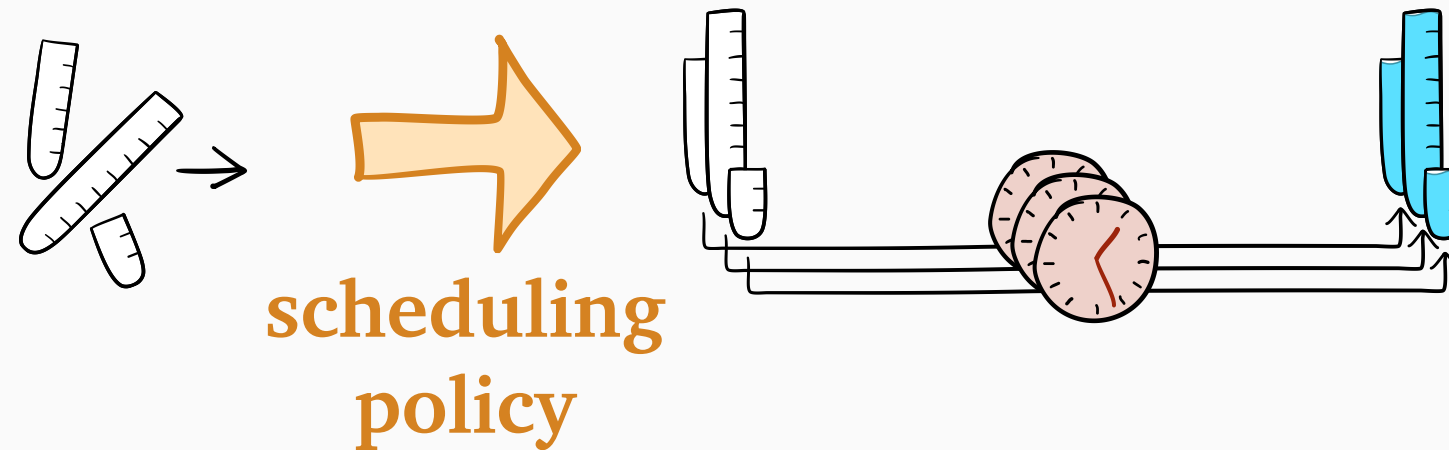
## Problem:

- Theory limited to “simple” **policies**
- In practice, “complex” **policies** matter



What makes a **policy** “simple” or “complex”?

# Goal: analyze impact of **policy**



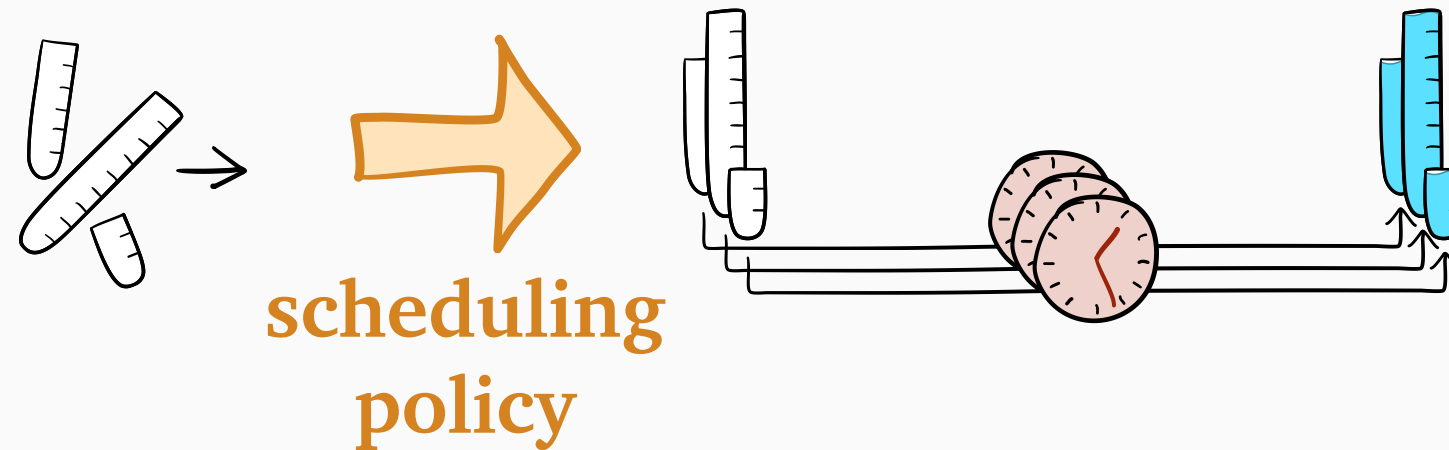
## Problem:

- Theory limited to “simple” **policies**
- In practice, “complex” **policies** matter

? What makes a **policy** “simple” or “complex”?

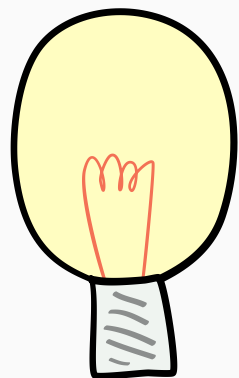
? What makes “complex” **policies** hard to analyze?

# Goal: analyze impact of **policy**



## Problem:

- Theory limited to “simple” **policies**
- In practice, “complex” **policies** matter



**SOAP** mindset:  
unifying *language* to  
describe **policies**

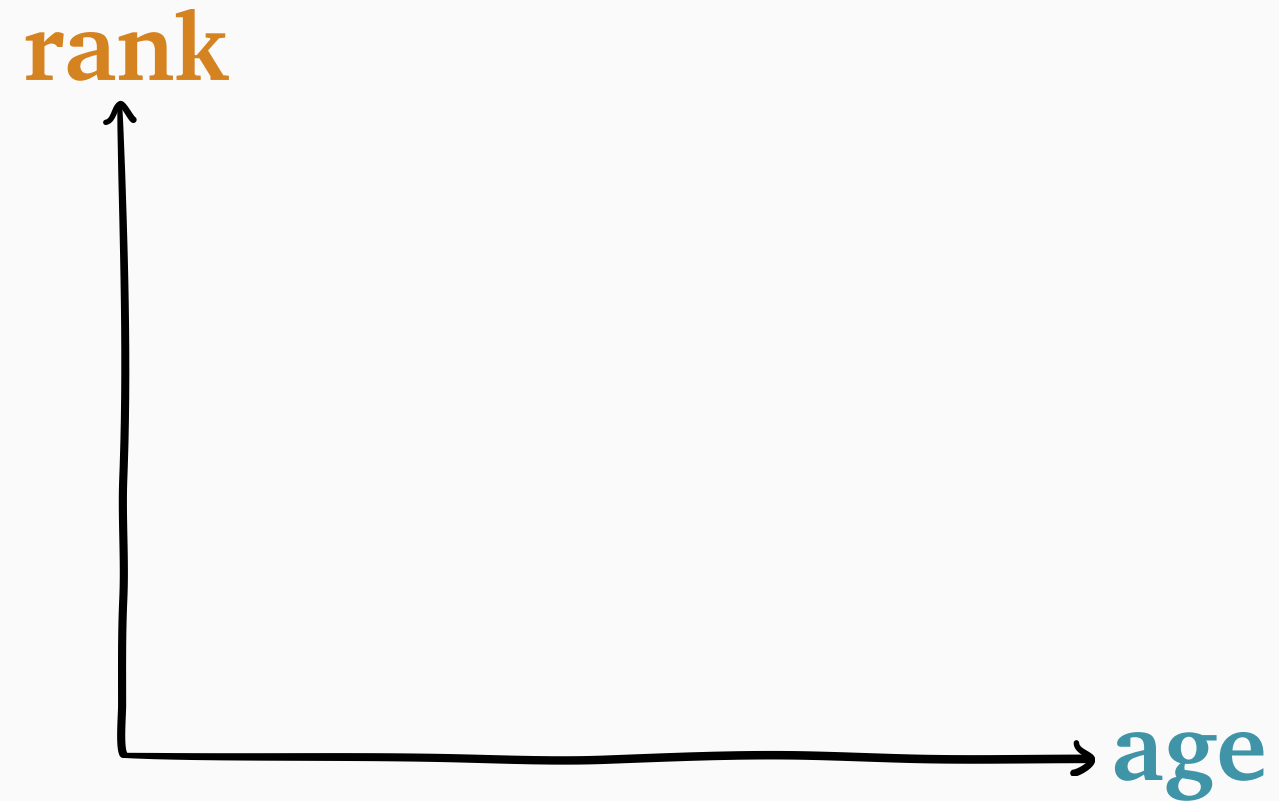


What makes a **policy**  
“simple” or “complex”?



What makes “complex”  
**policies** hard to analyze?

# Unifying language: **rank** functions

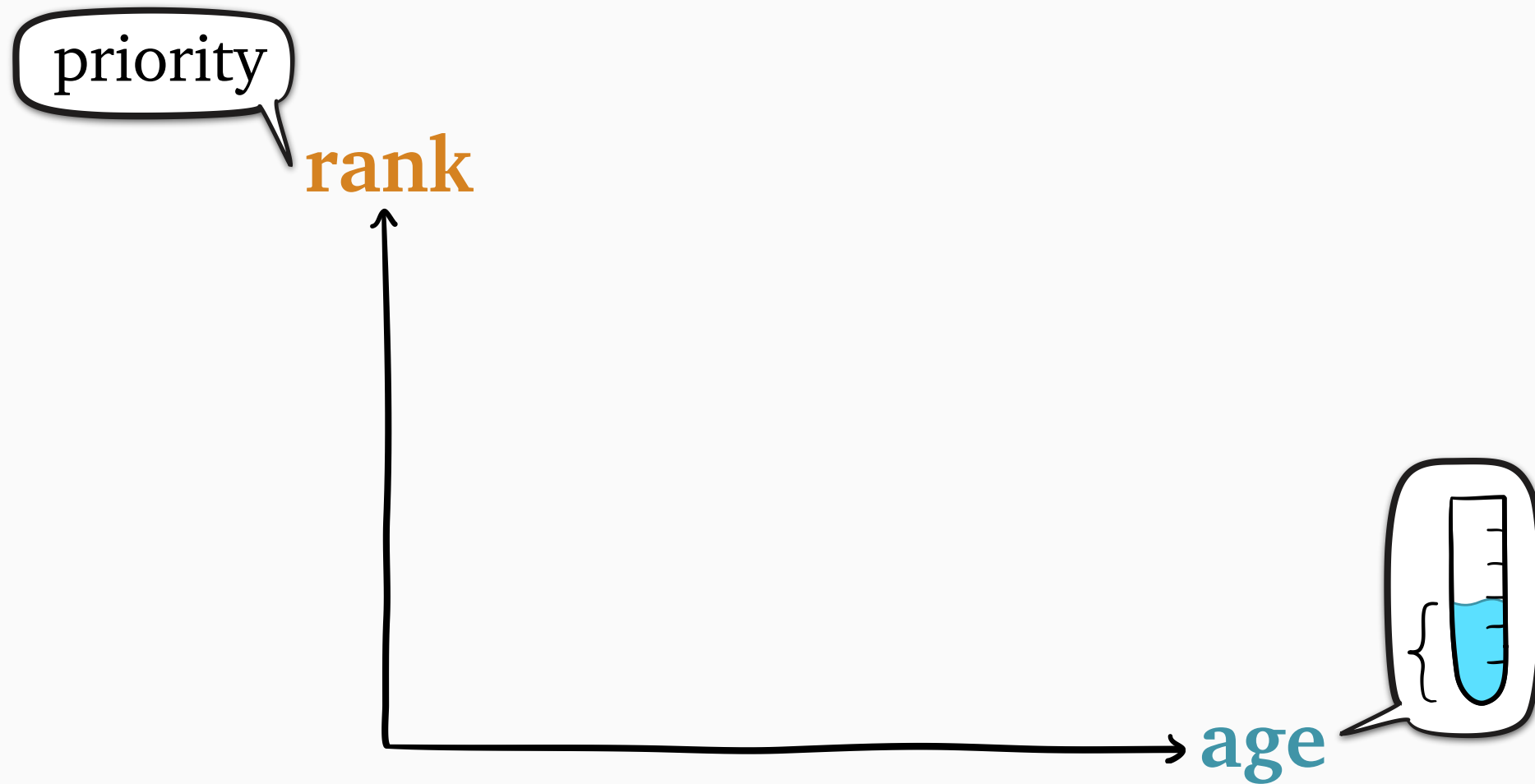




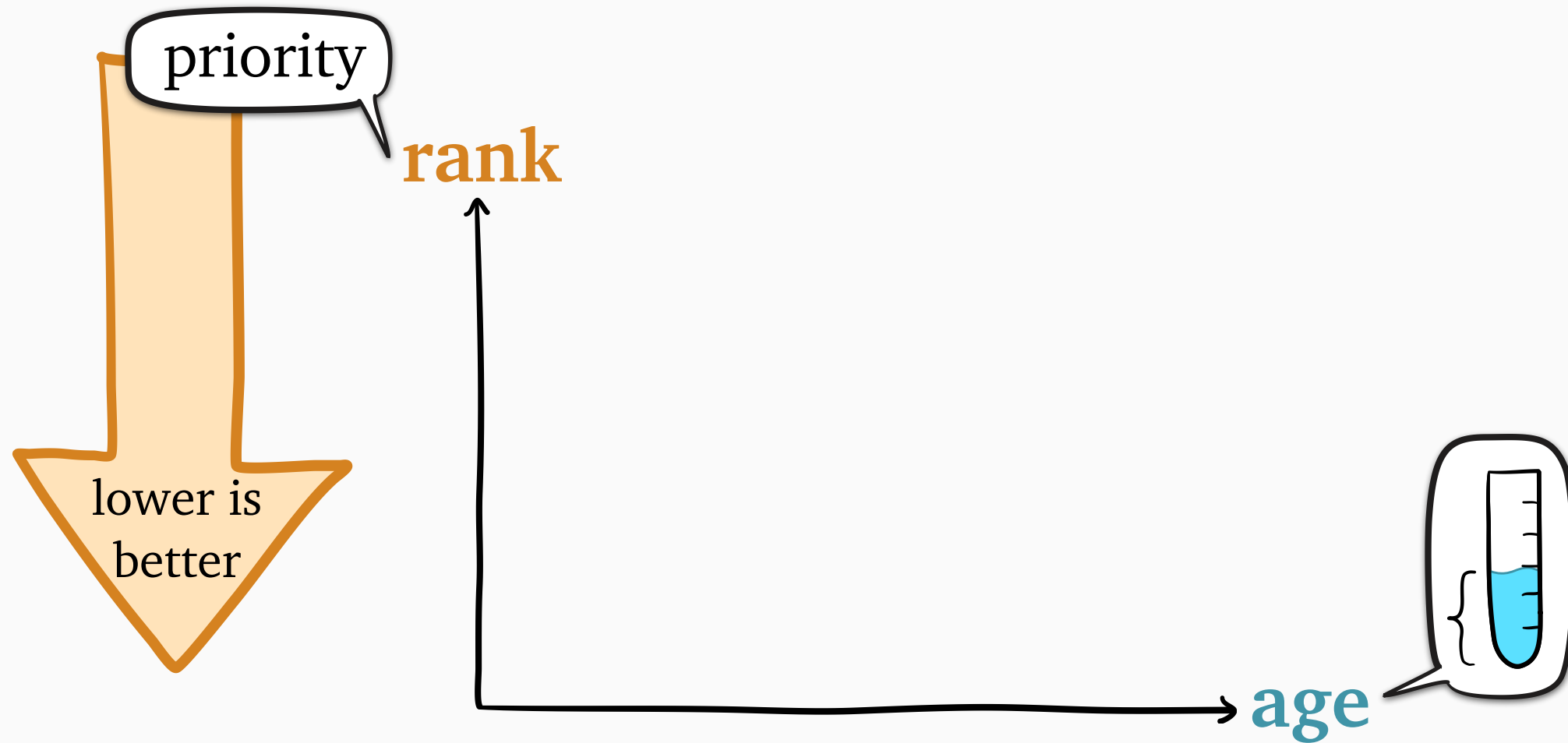
# Unifying language: **rank** functions



# Unifying language: **rank** functions



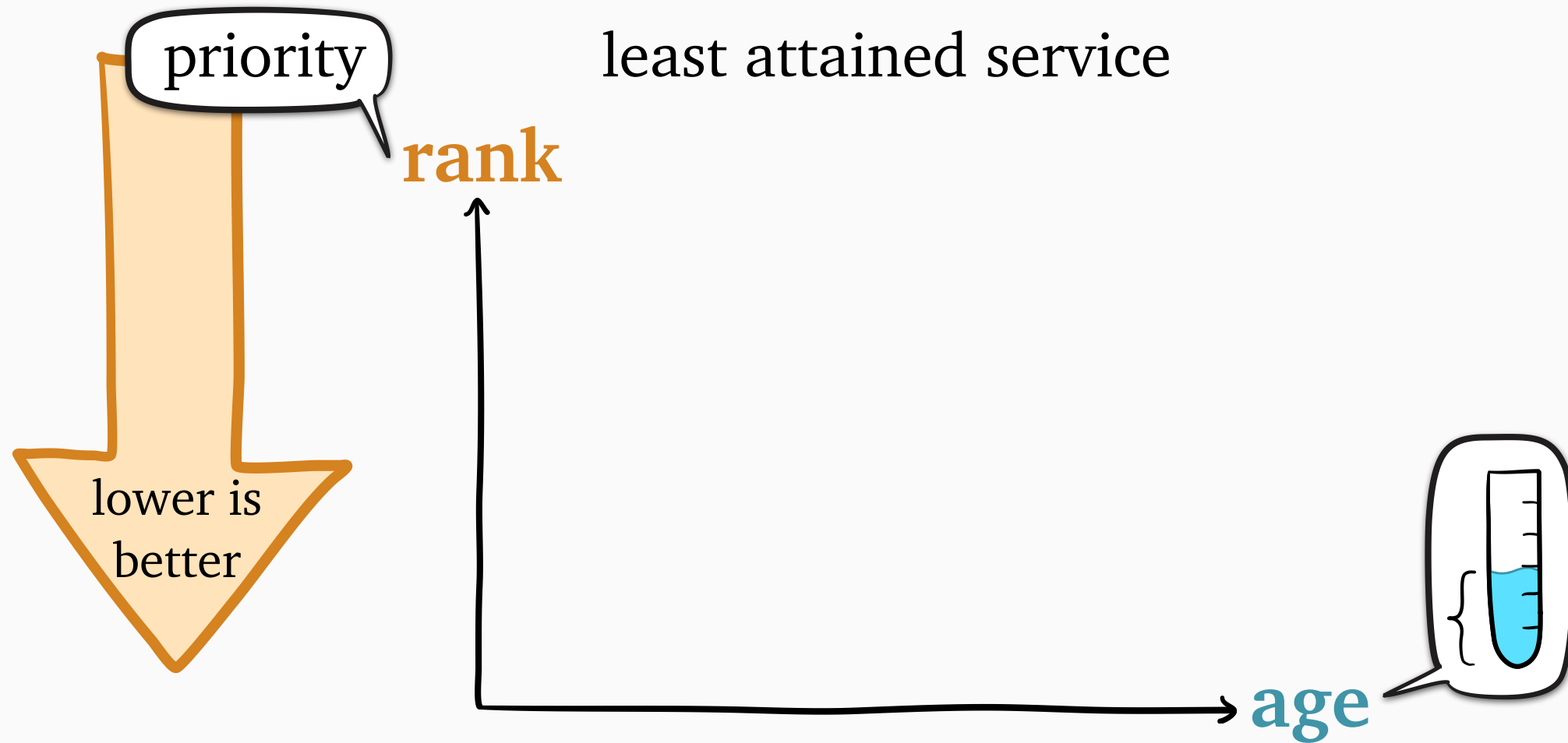
# Unifying language: **rank** functions



# Unifying language: **rank** functions

**LAS**

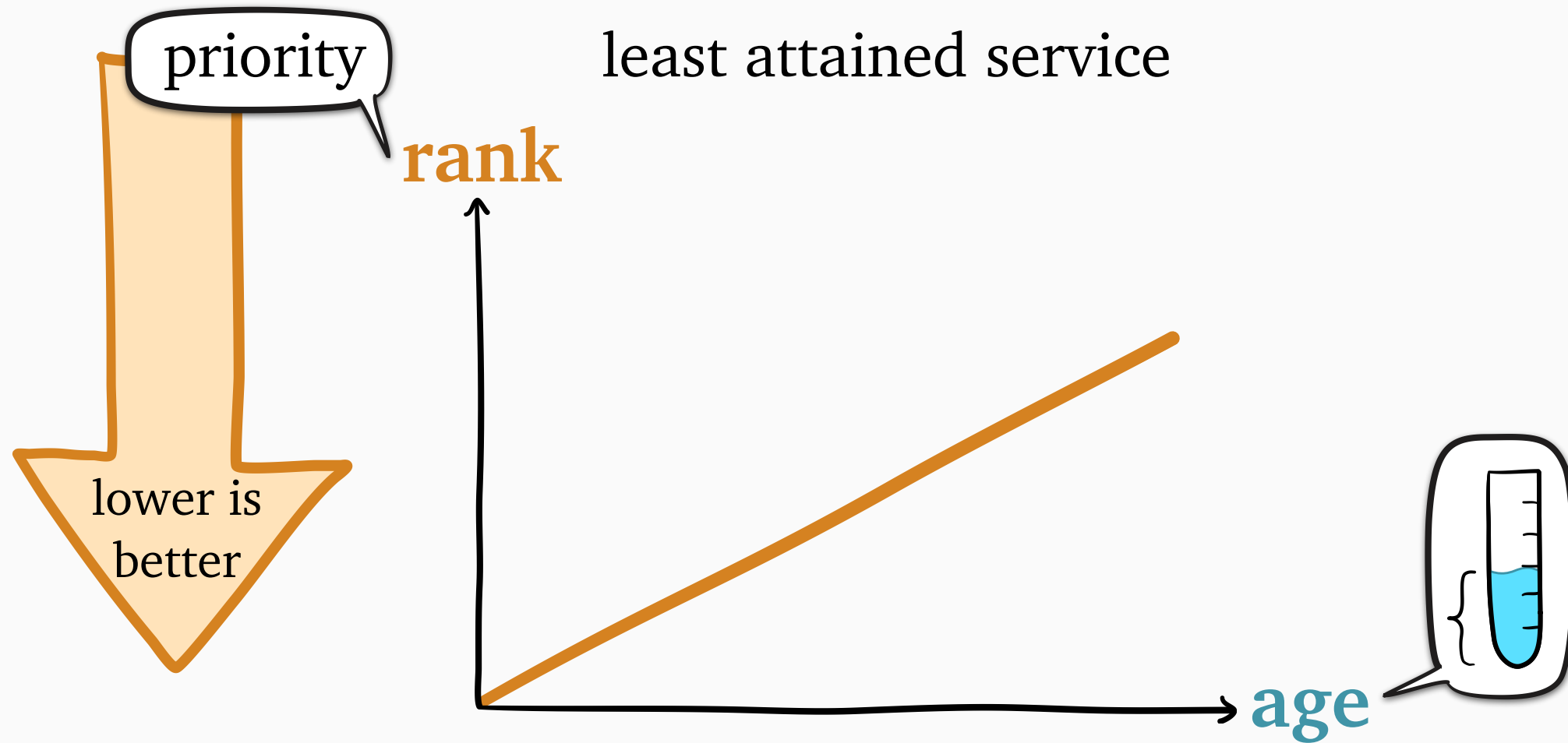
least attained service



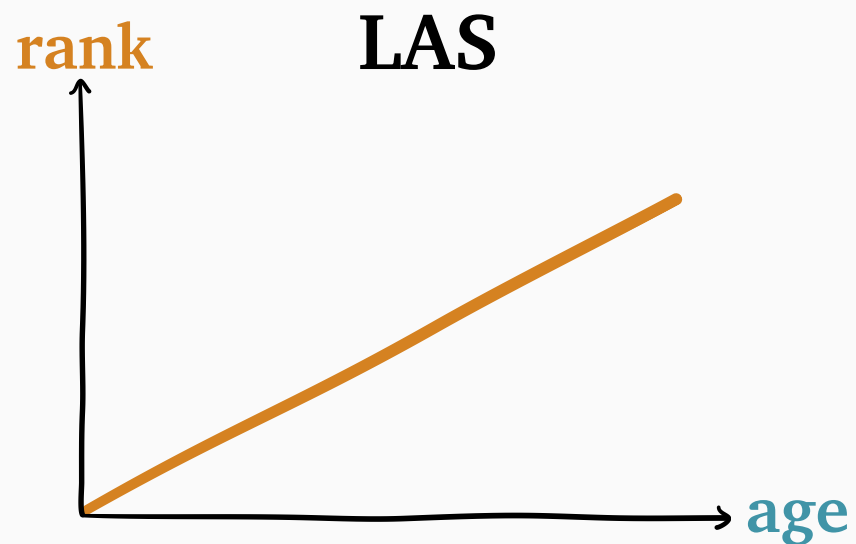
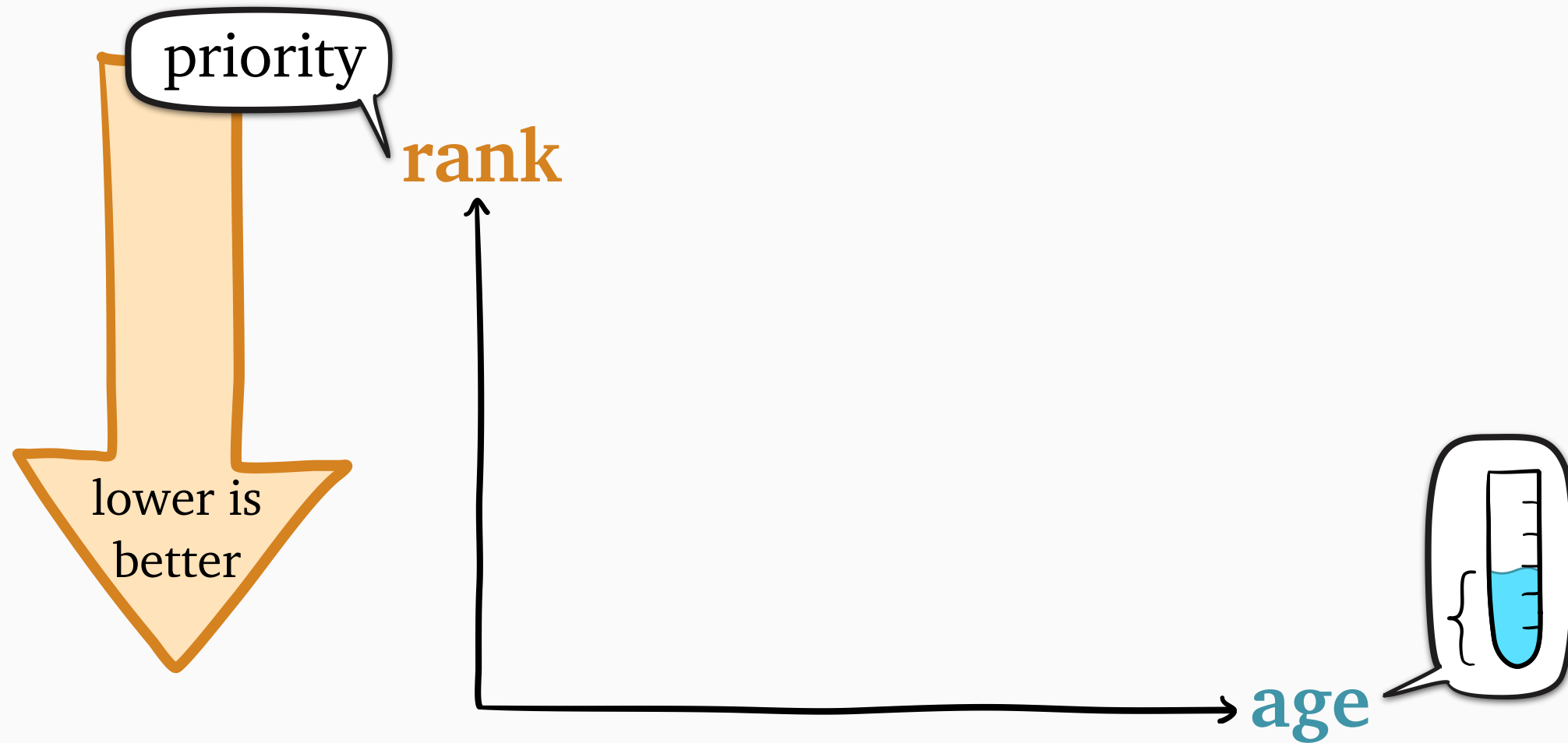
# Unifying language: **rank** functions

**LAS**

least attained service

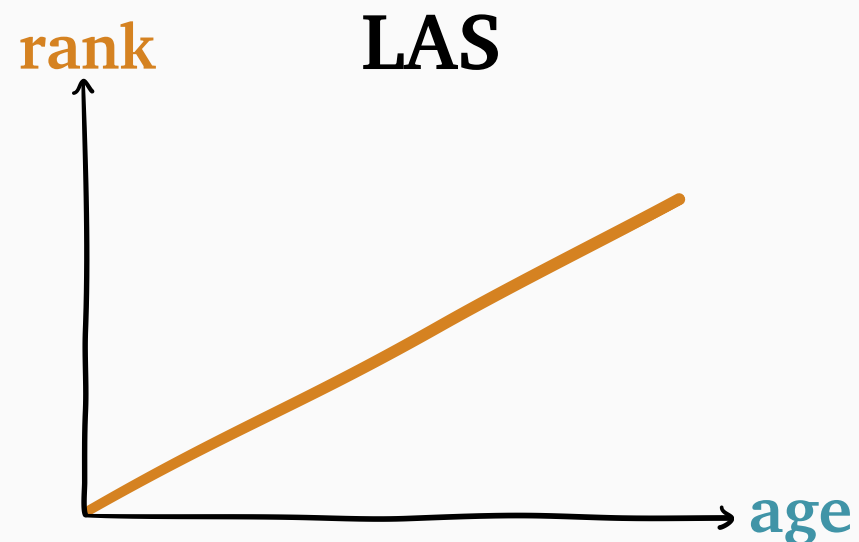
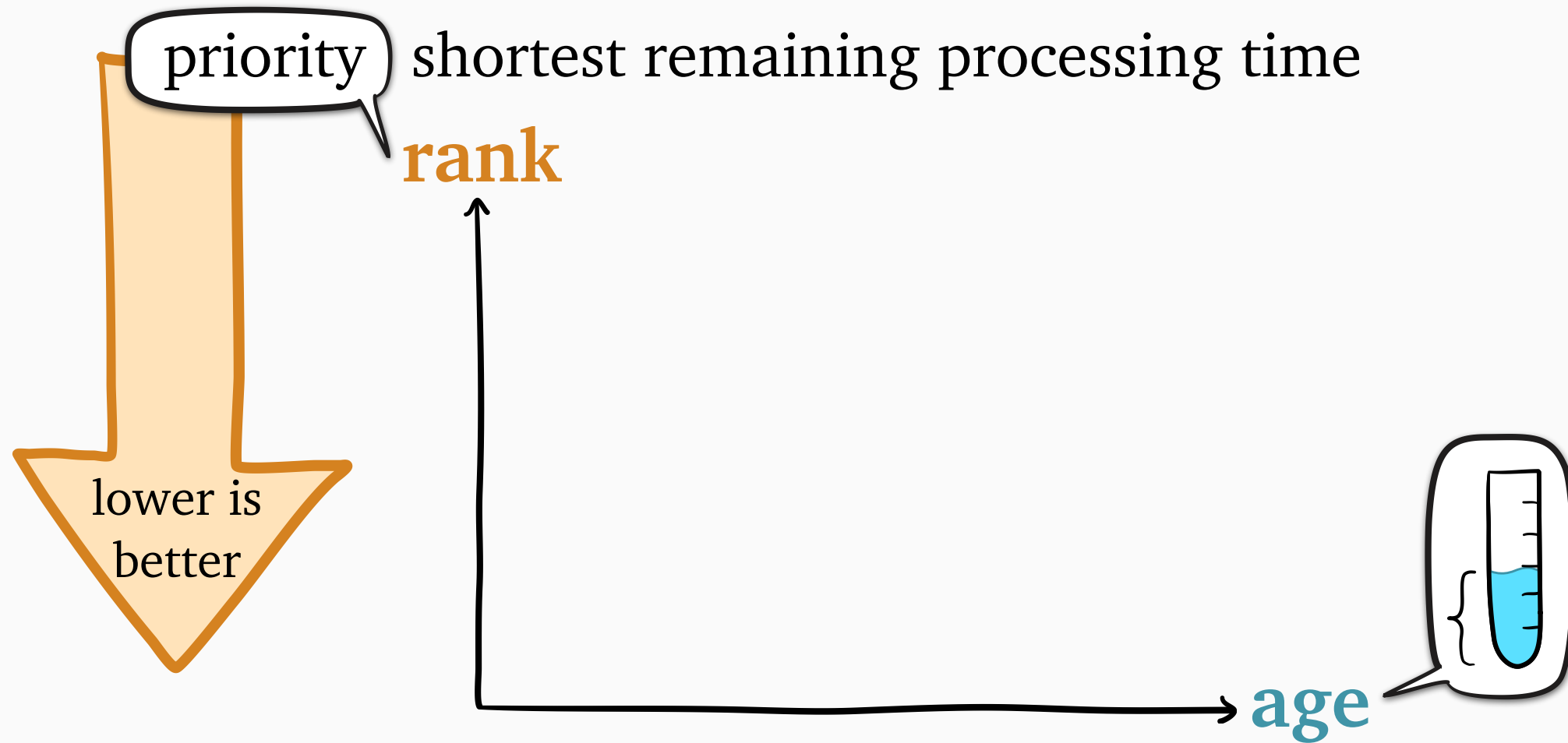


# Unifying language: **rank** functions



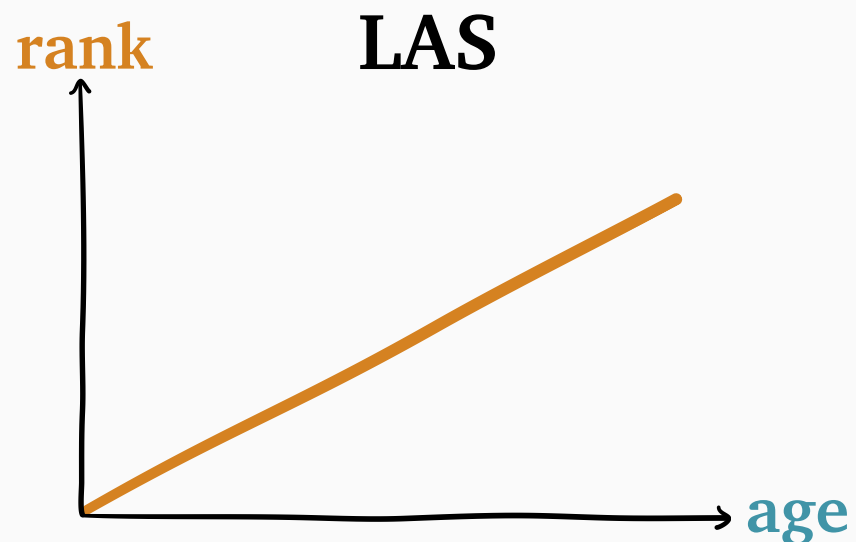
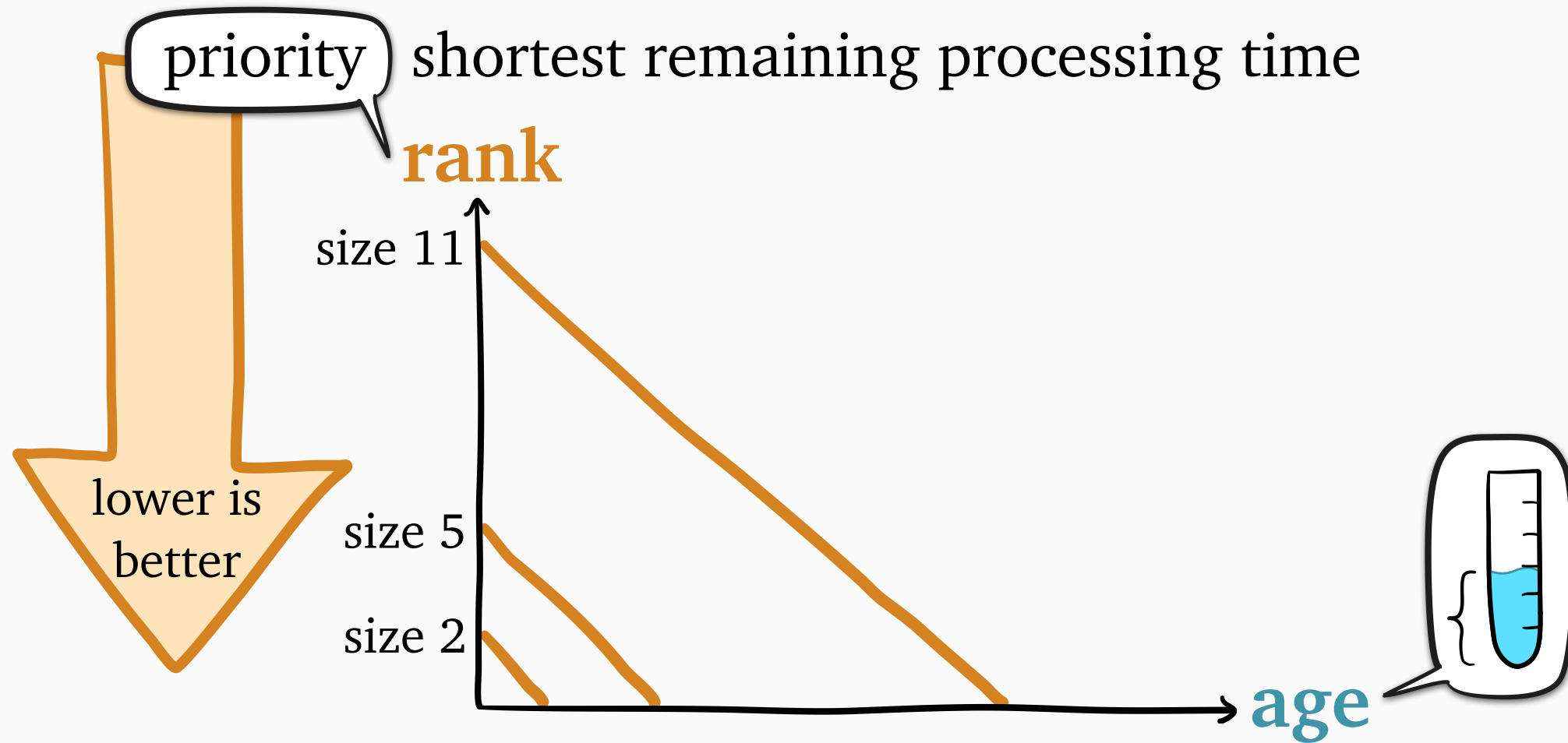
# Unifying language: **rank** functions

## SRPT



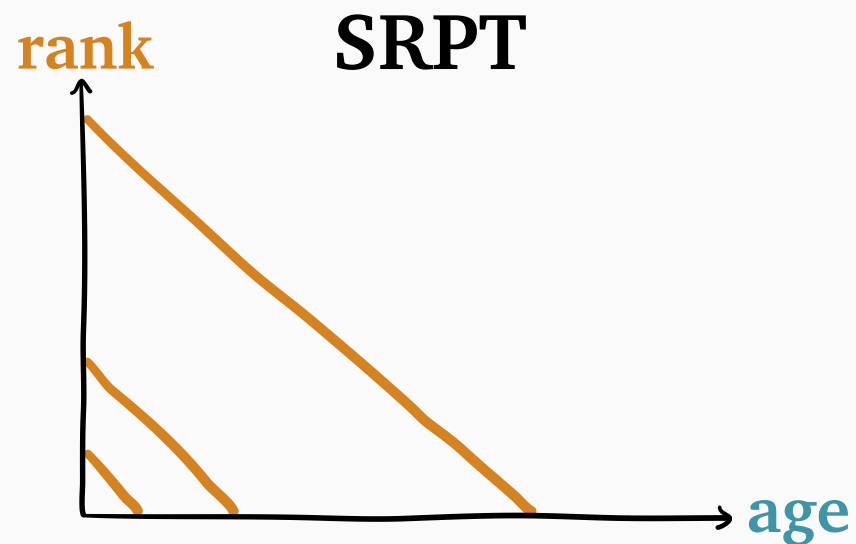
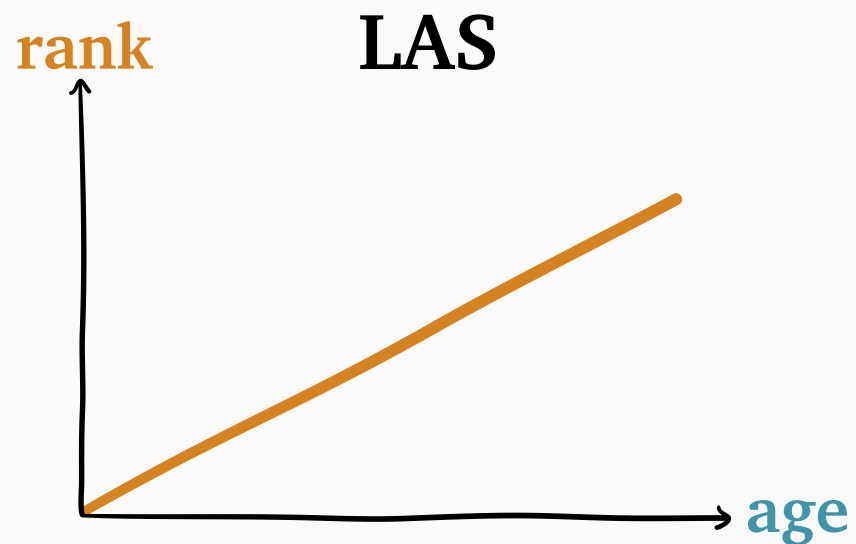
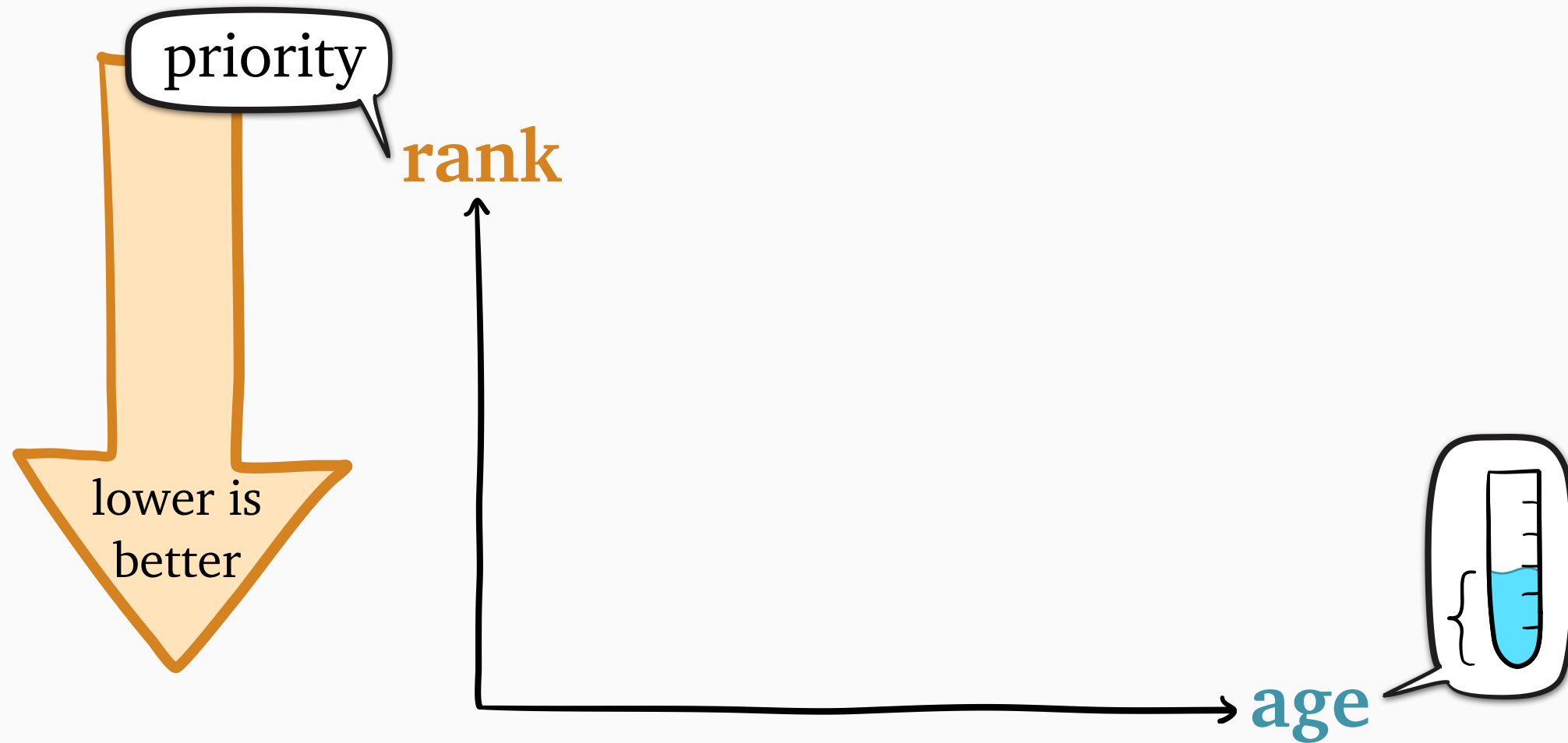
# Unifying language: **rank** functions

## SRPT





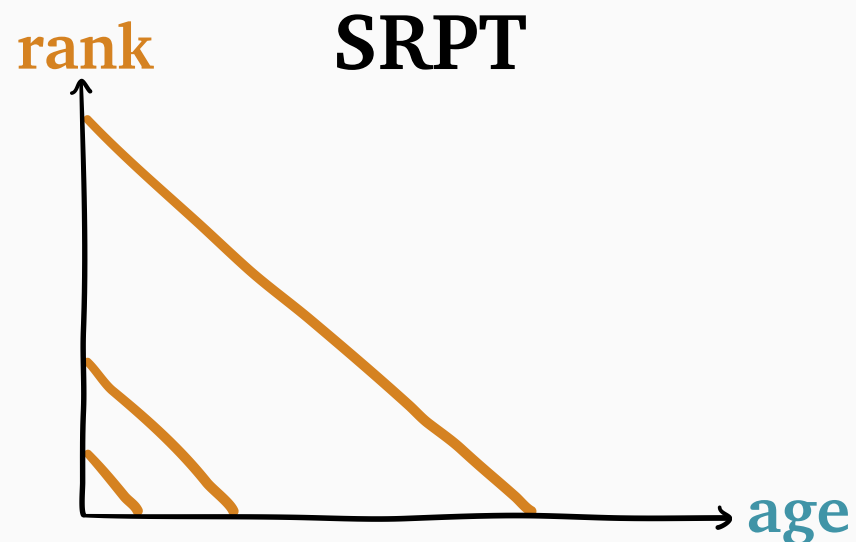
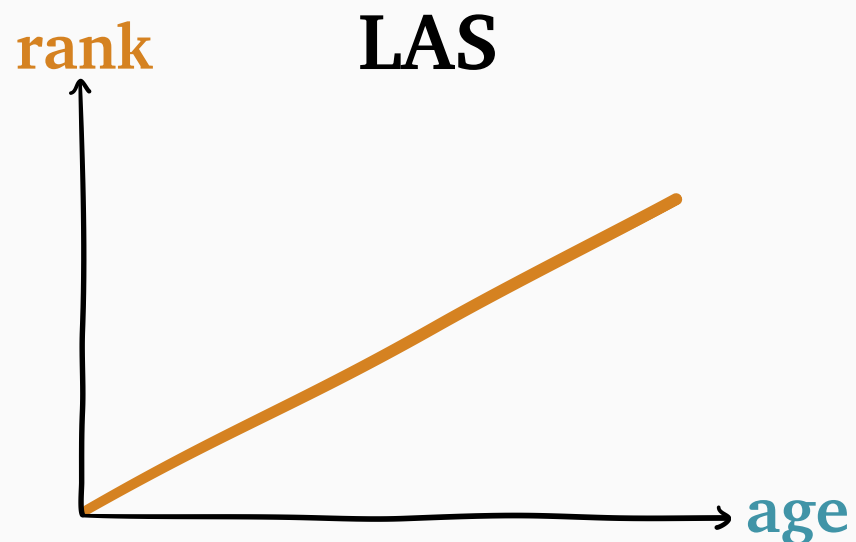
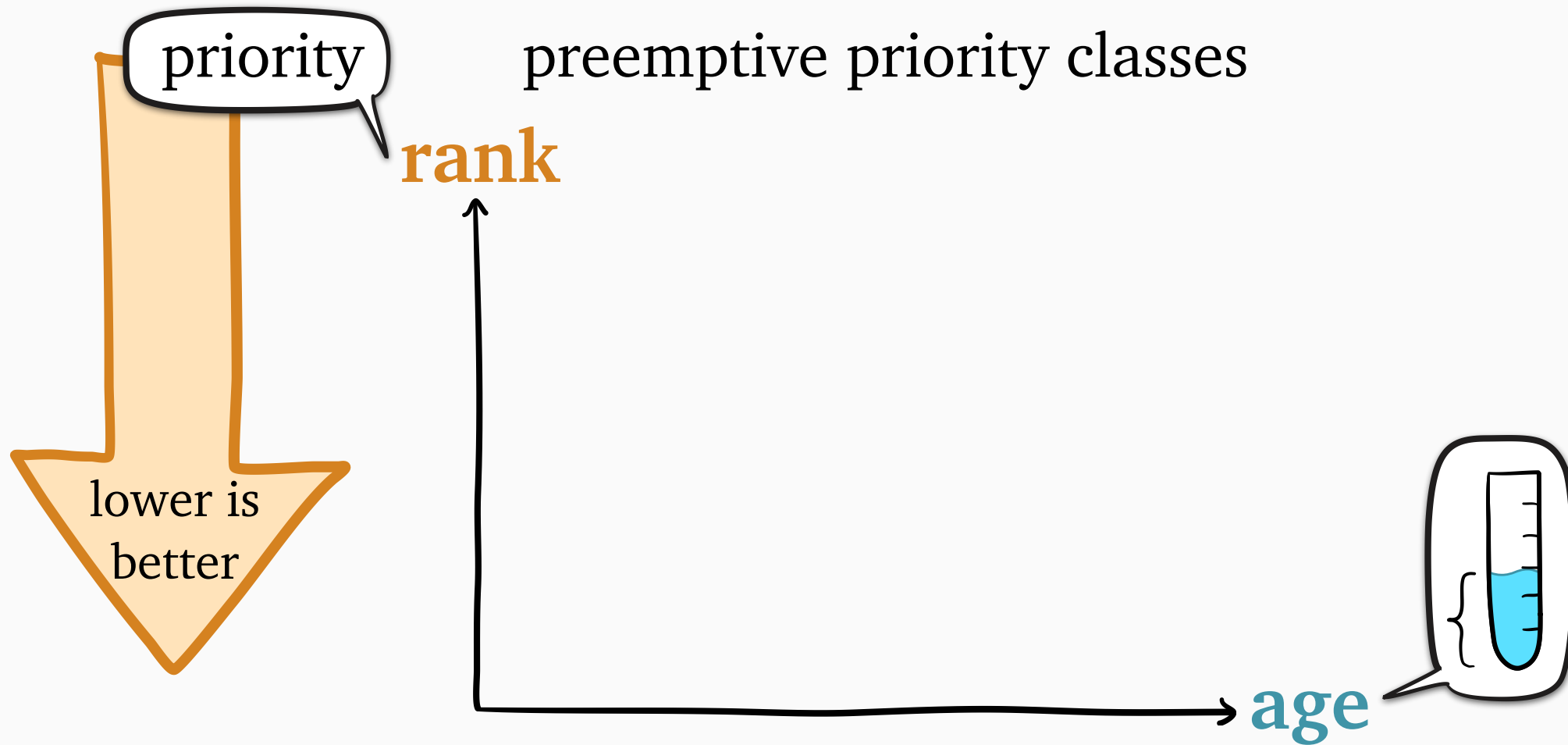
# Unifying language: **rank** functions



# Unifying language: **rank** functions

## P-Prio

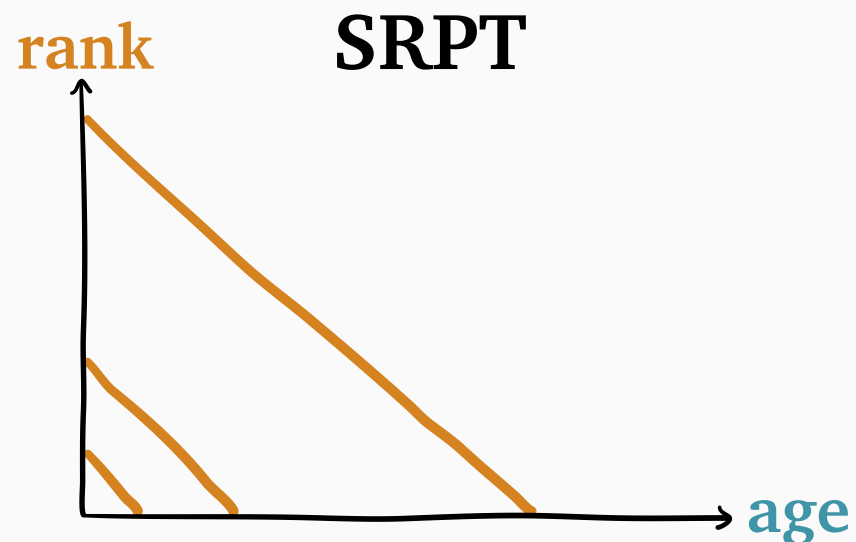
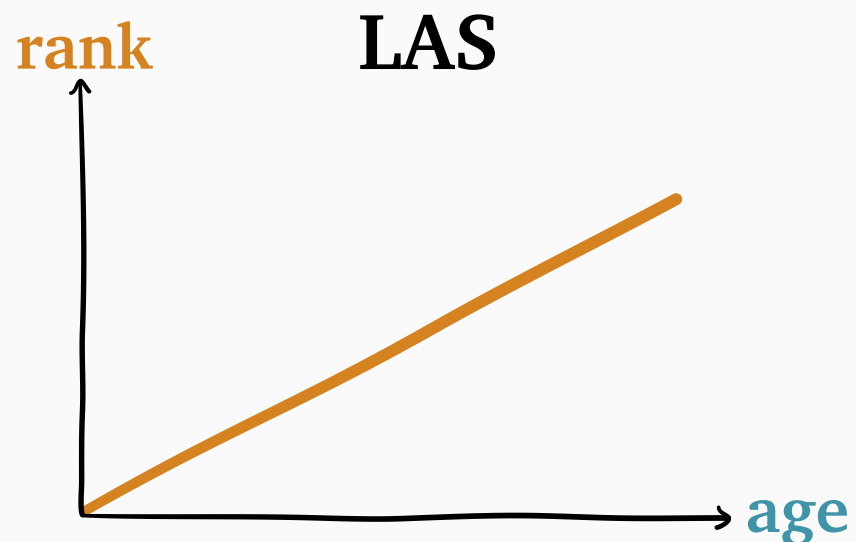
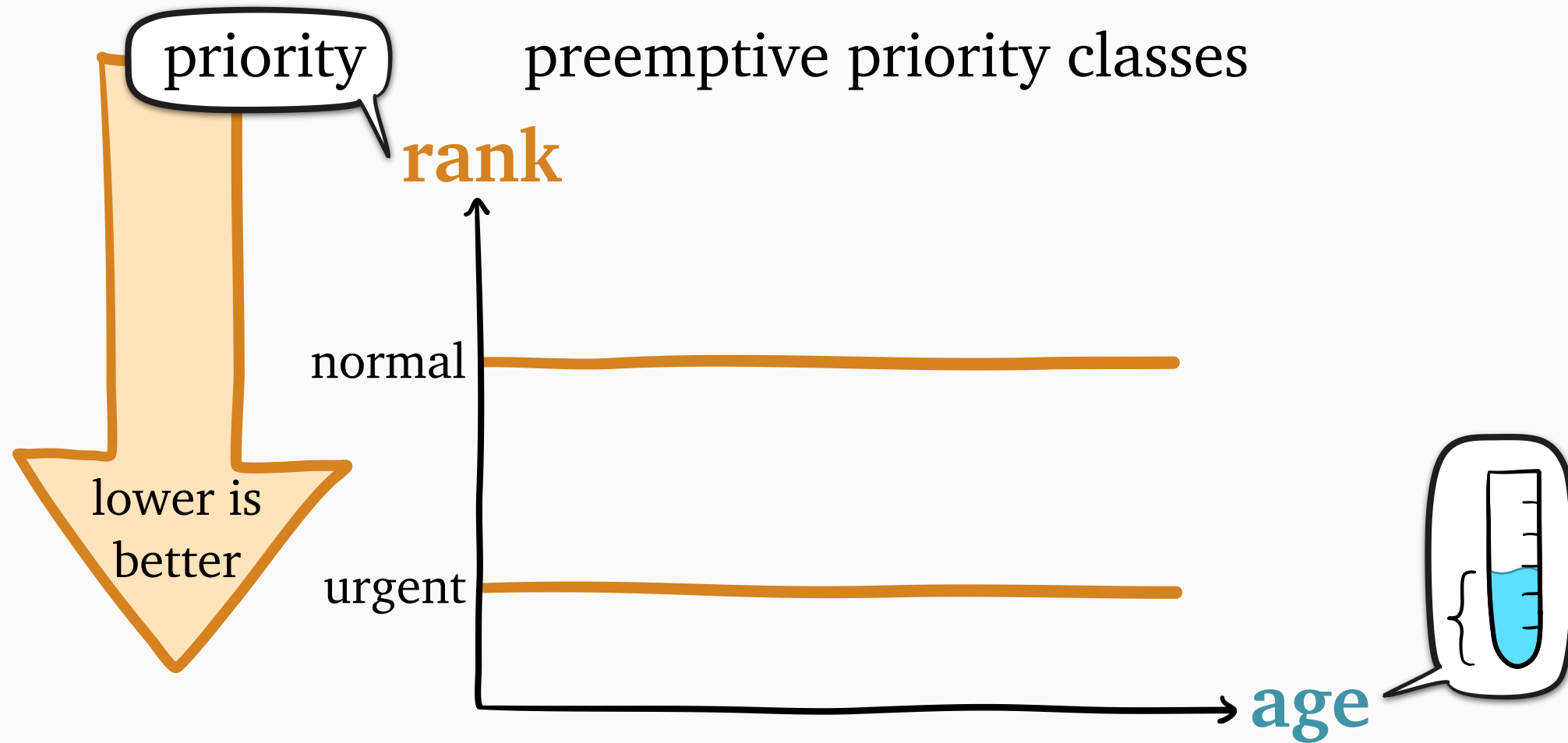
preemptive priority classes



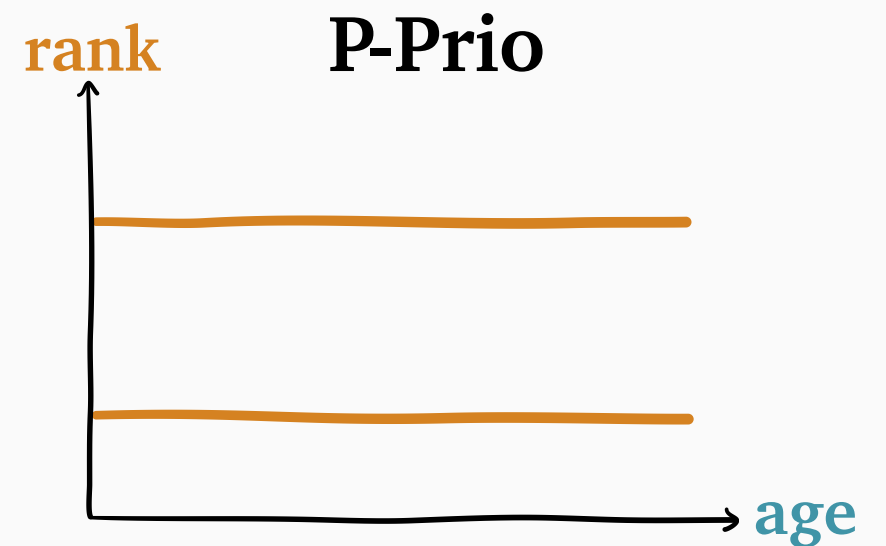
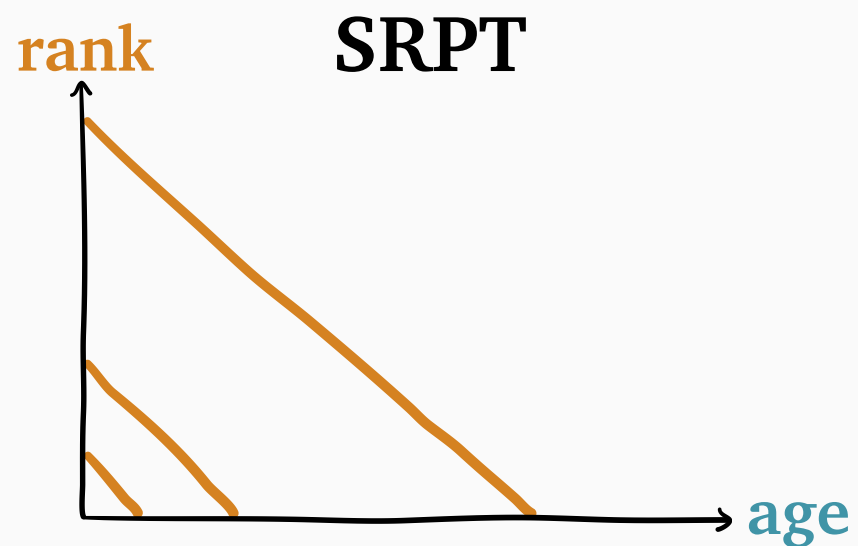
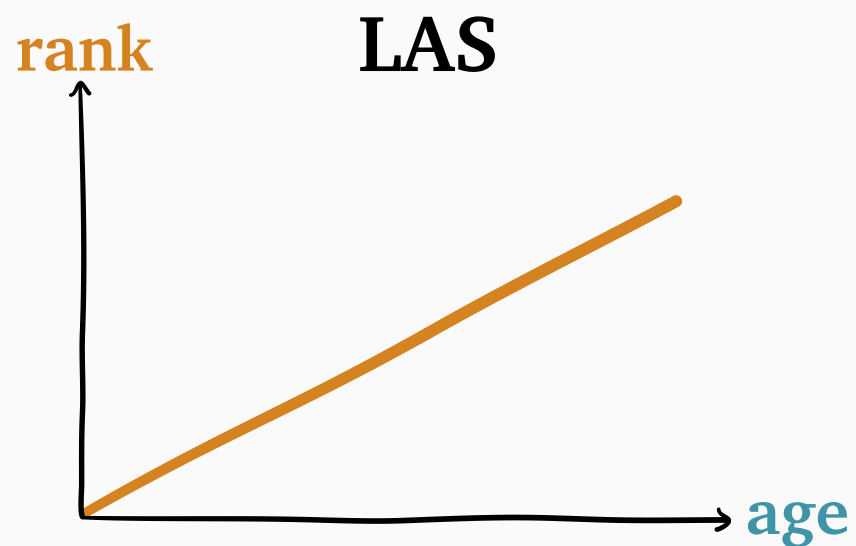
# Unifying language: **rank** functions

## P-Prio

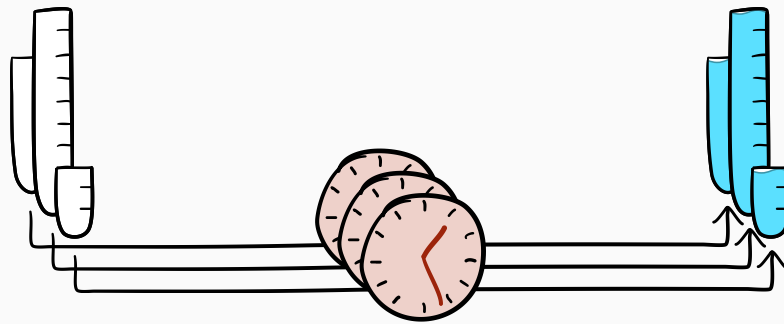
preemptive priority classes



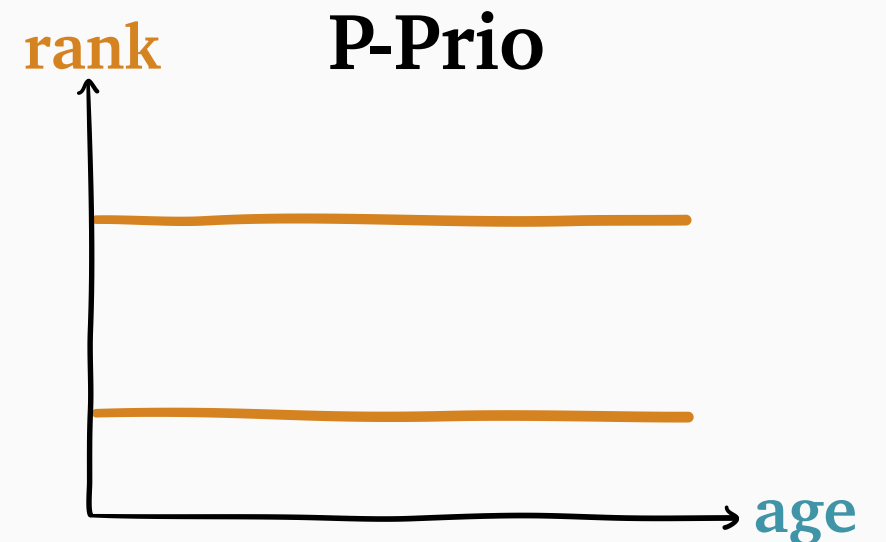
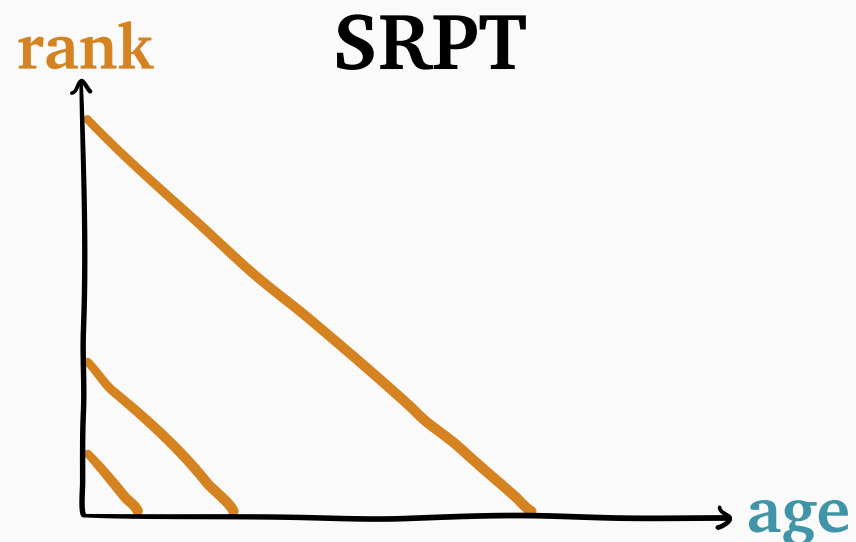
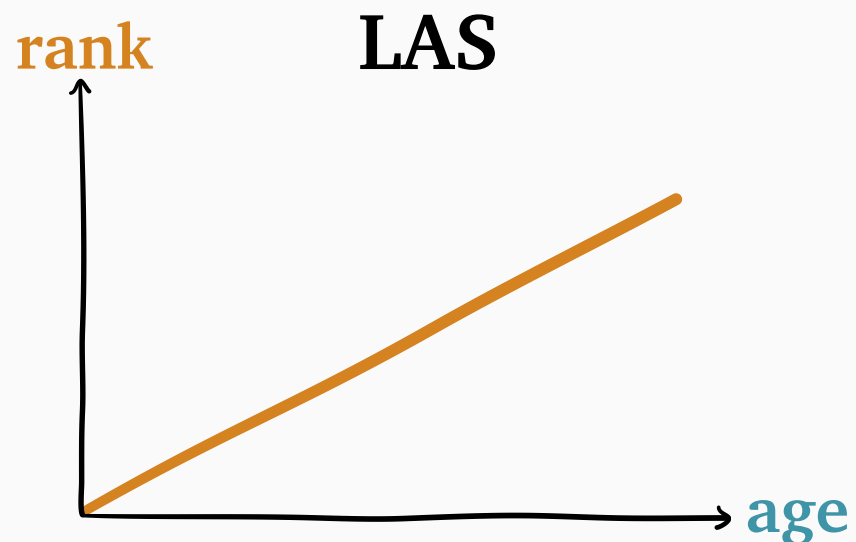
# Unifying language: **rank** functions



# Unifying language: **rank** functions

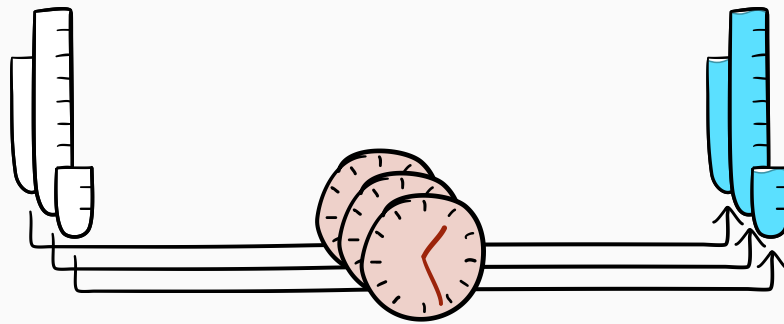


response time  $T$  known

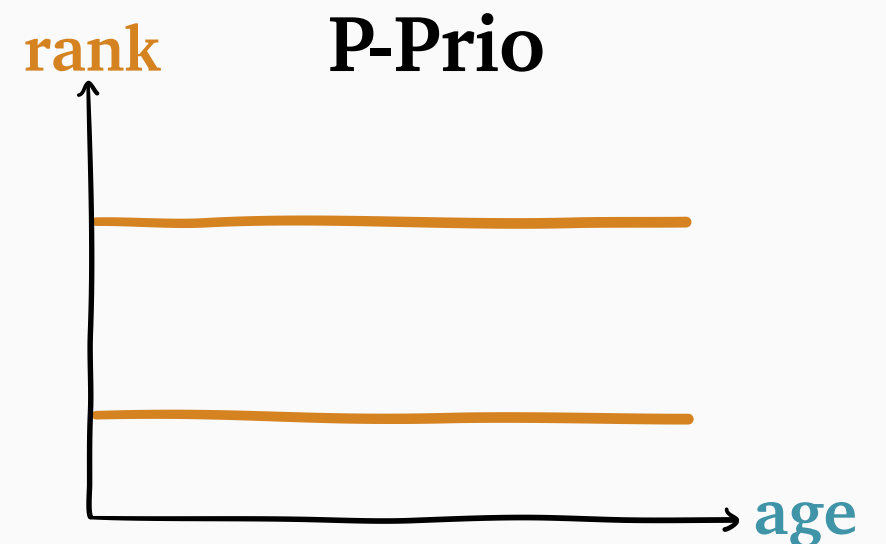
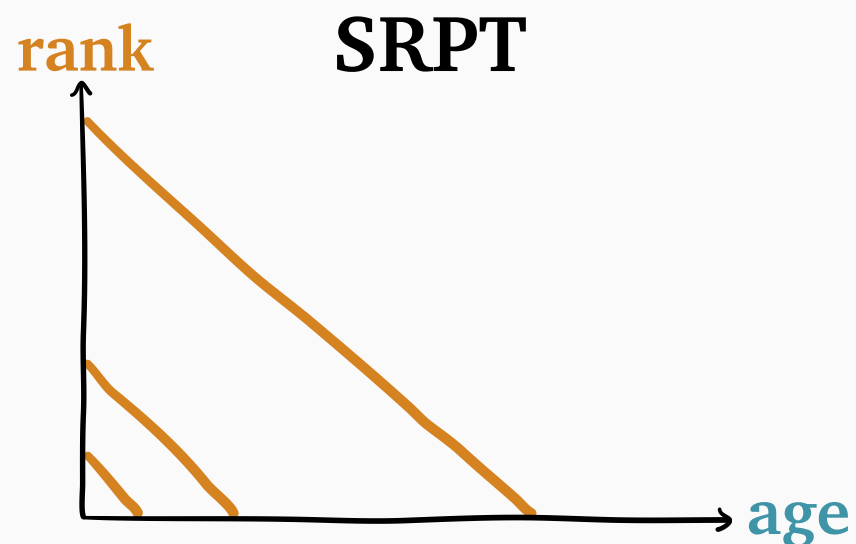
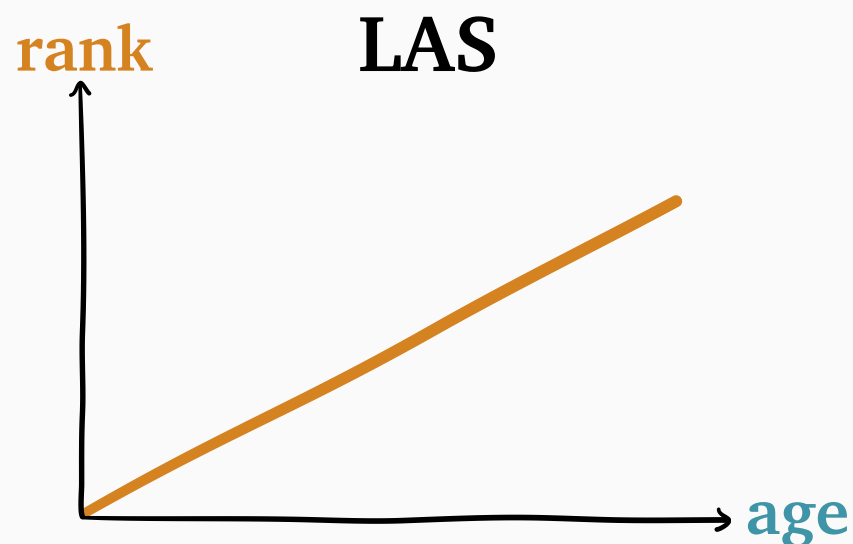


# Unifying language: **rank** functions

“Simple” **policy** =                     ? **rank** function

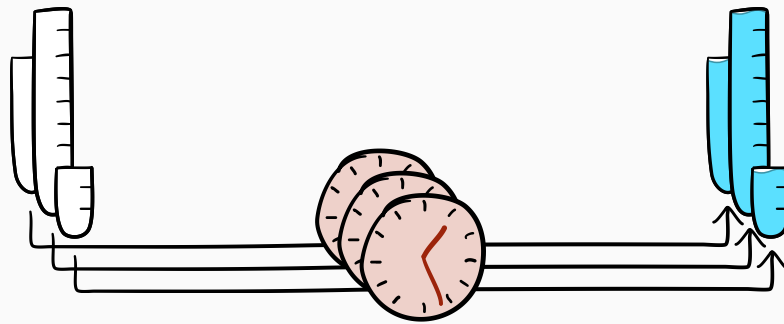


response time  $T$  known

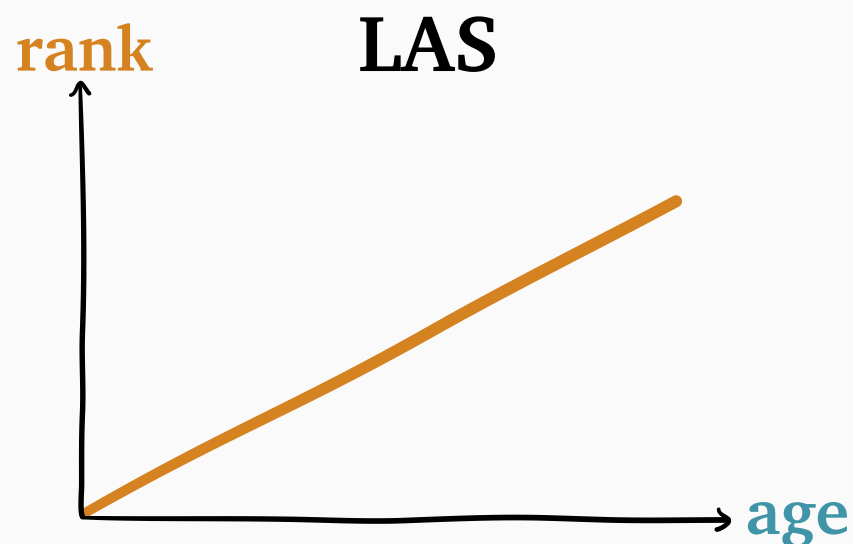


# Unifying language: **rank** functions

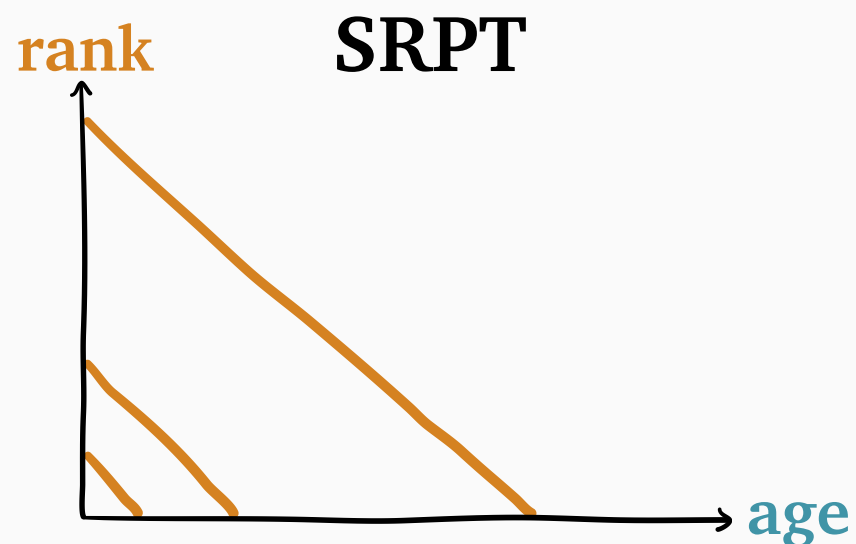
“Simple” **policy** = monotonic **rank** function



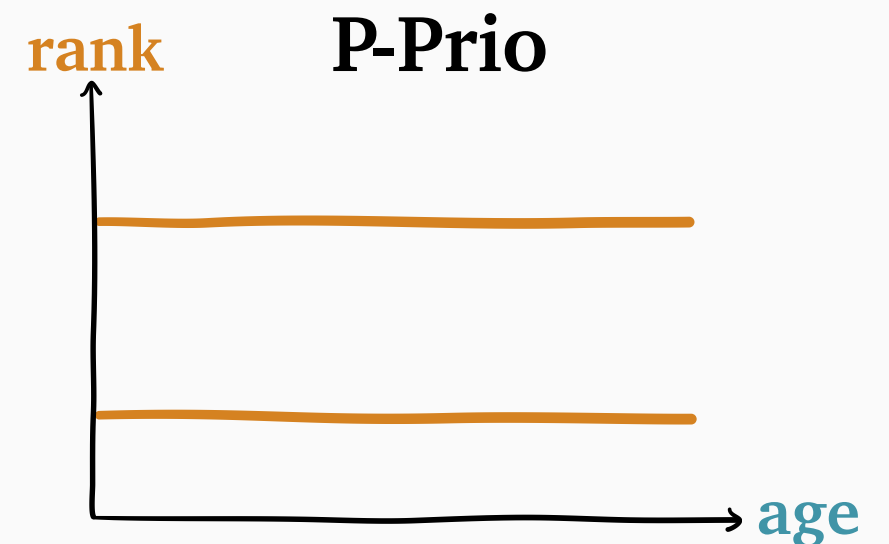
response time  $T$  known



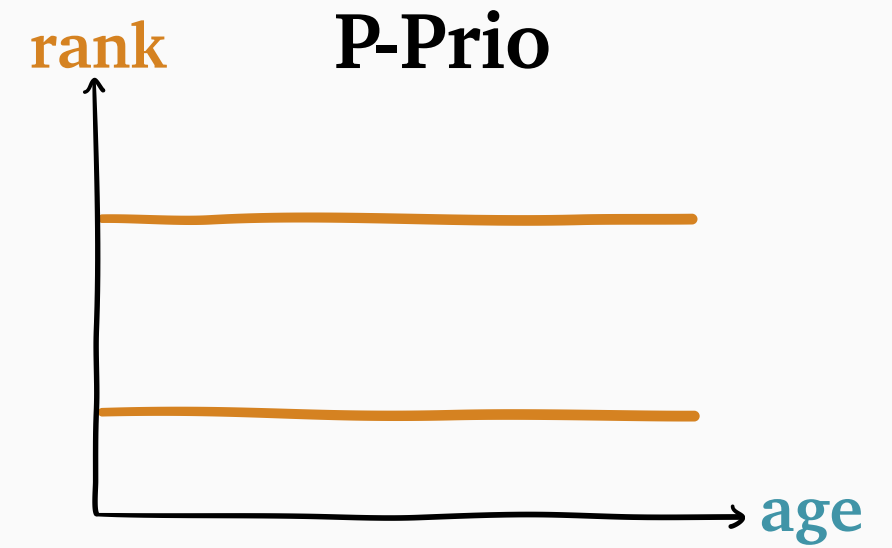
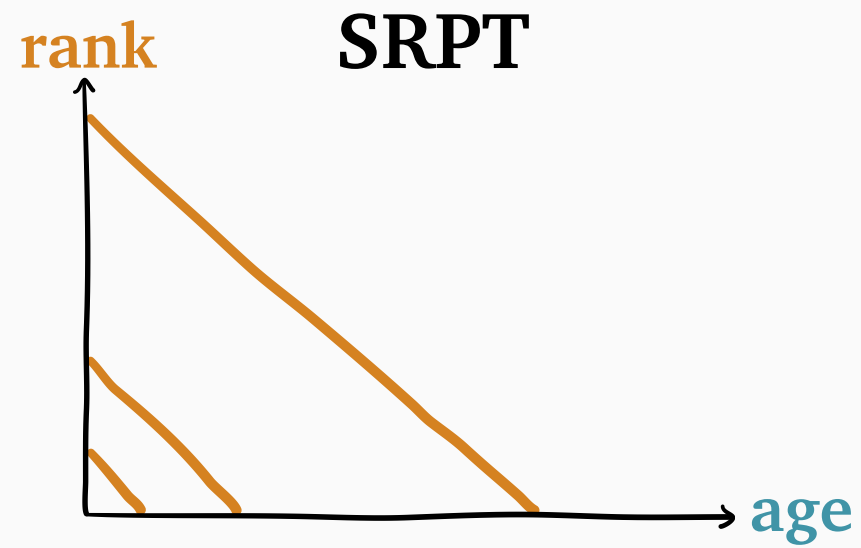
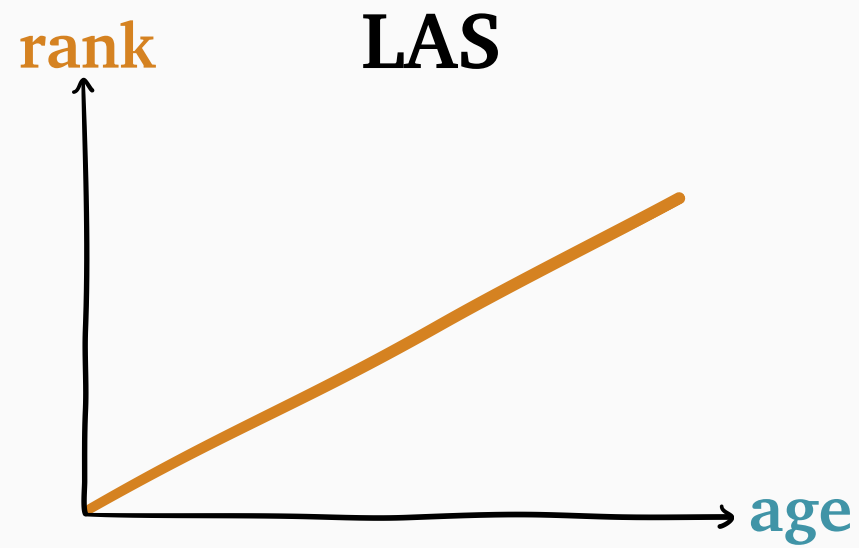
LAS



SRPT

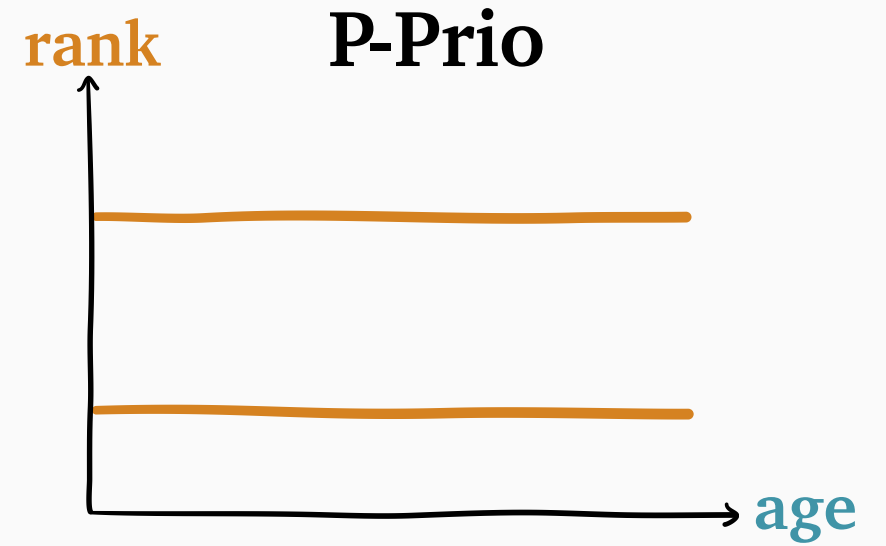
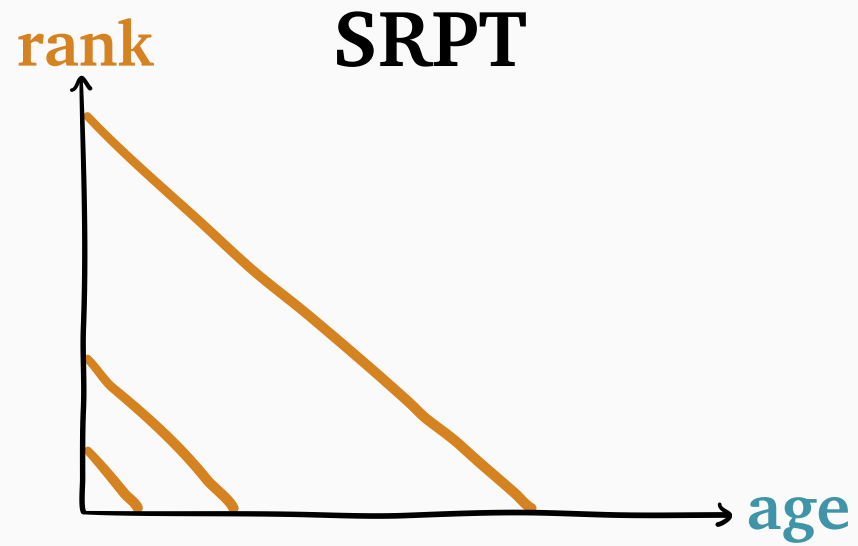
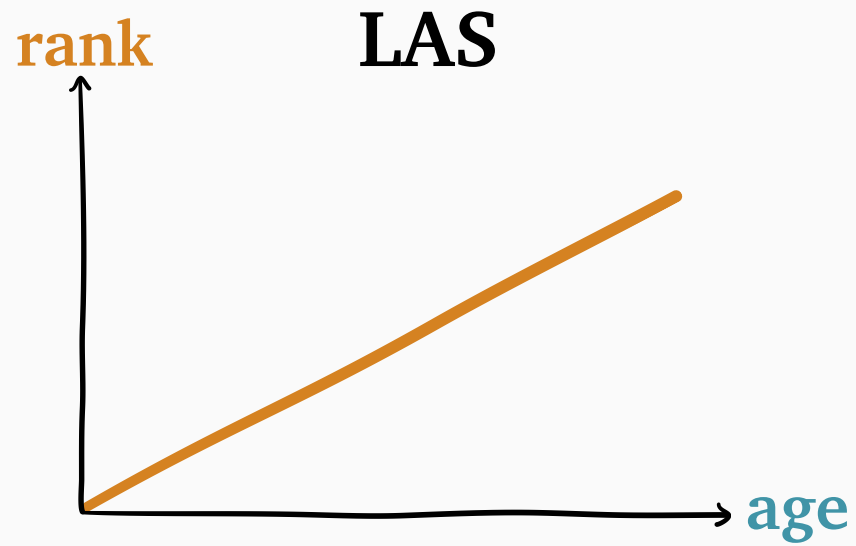


P-Prio



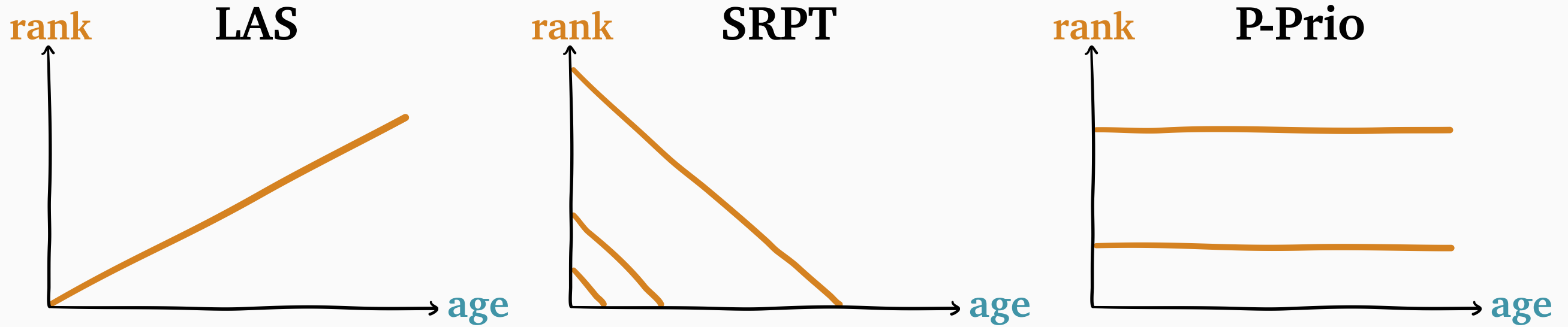
“Simple” = *monotonic* rank function





“Simple” = *monotonic* rank function

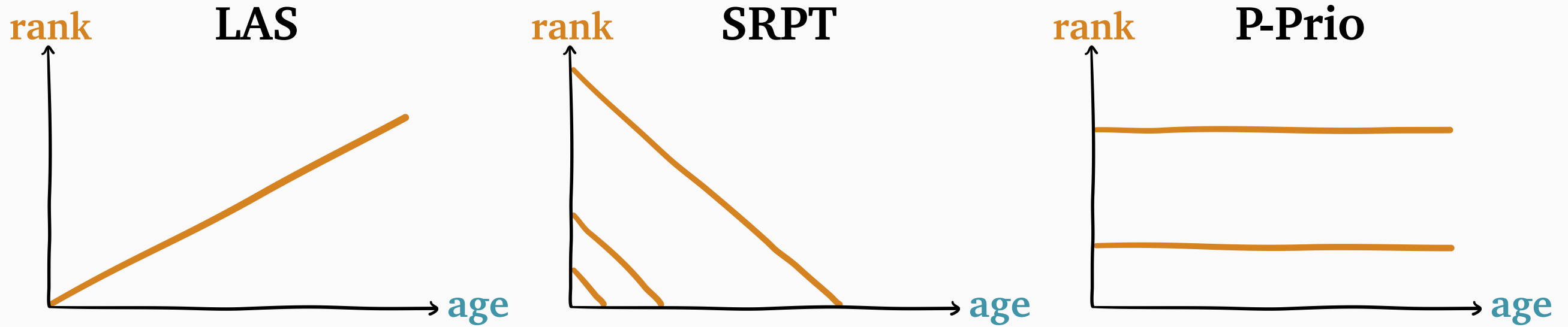
tractable before SOAP



“Simple” = *monotonic* **rank** function

tractable before **SOAP**

“Complex” = *nonmonotonic* **rank** function



“Simple” = *monotonic* **rank** function

tractable before **SOAP**

“Complex” = *nonmonotonic* **rank** function

*intractable* before **SOAP**

“Complex” = *nonmonotonic* **rank**

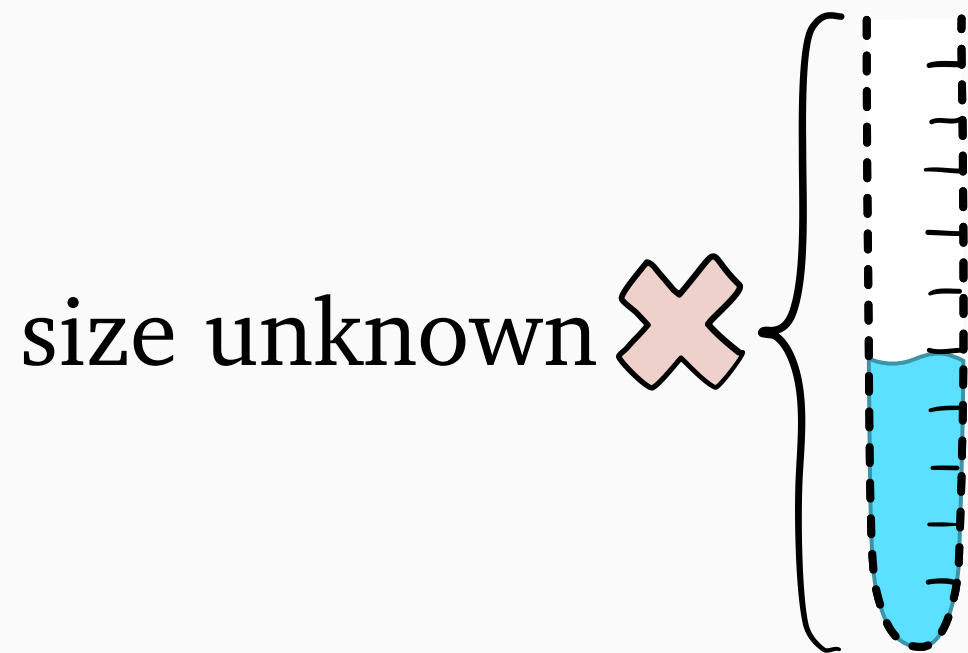
“Complex” = *nonmonotonic* **rank**



**Motivation:**  
unknown job sizes

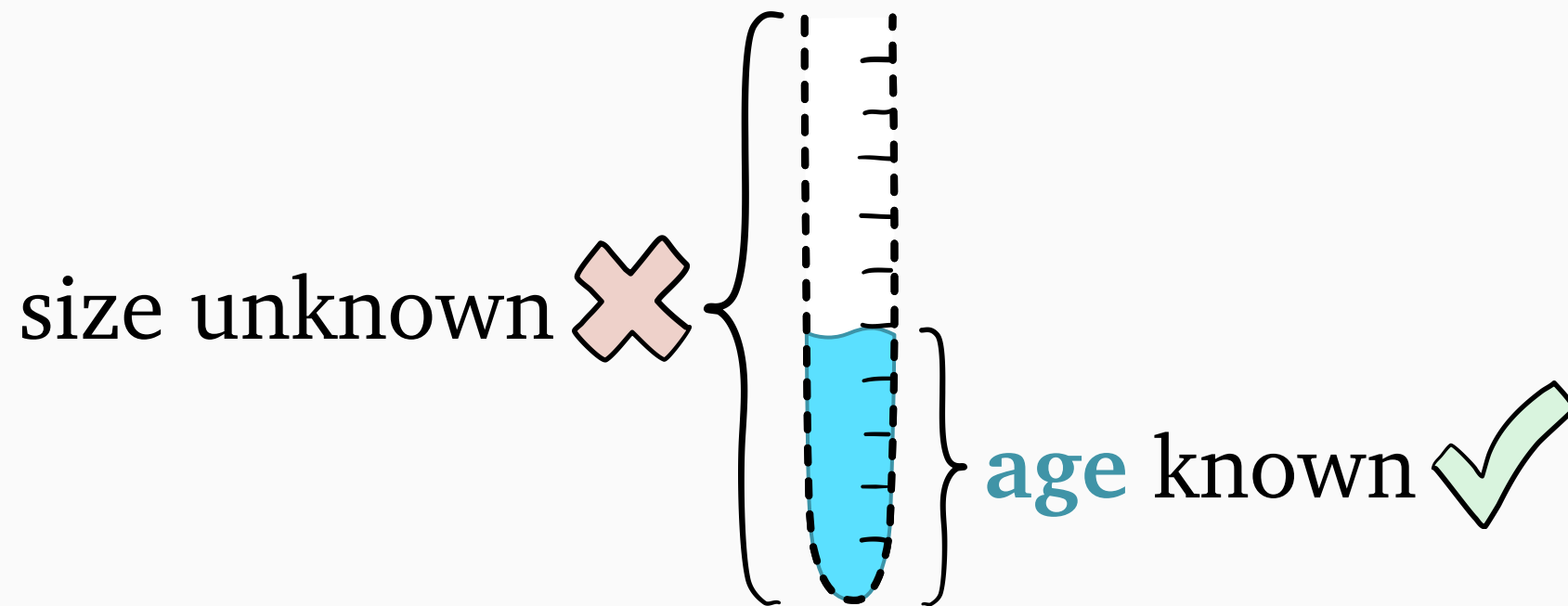
“Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes



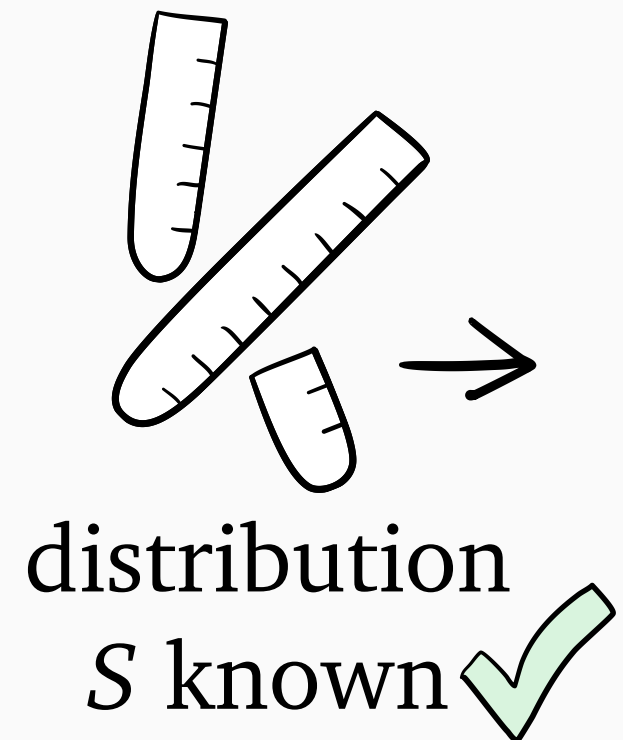
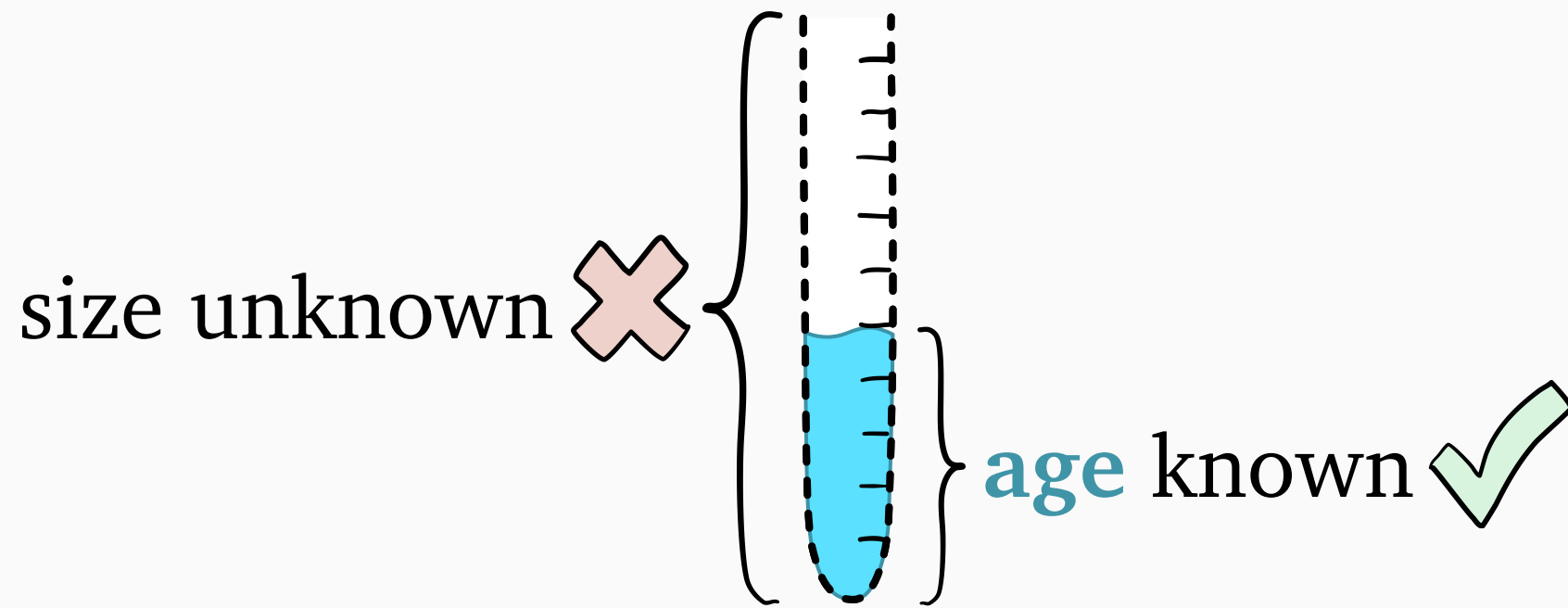
# “Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes



# “Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes



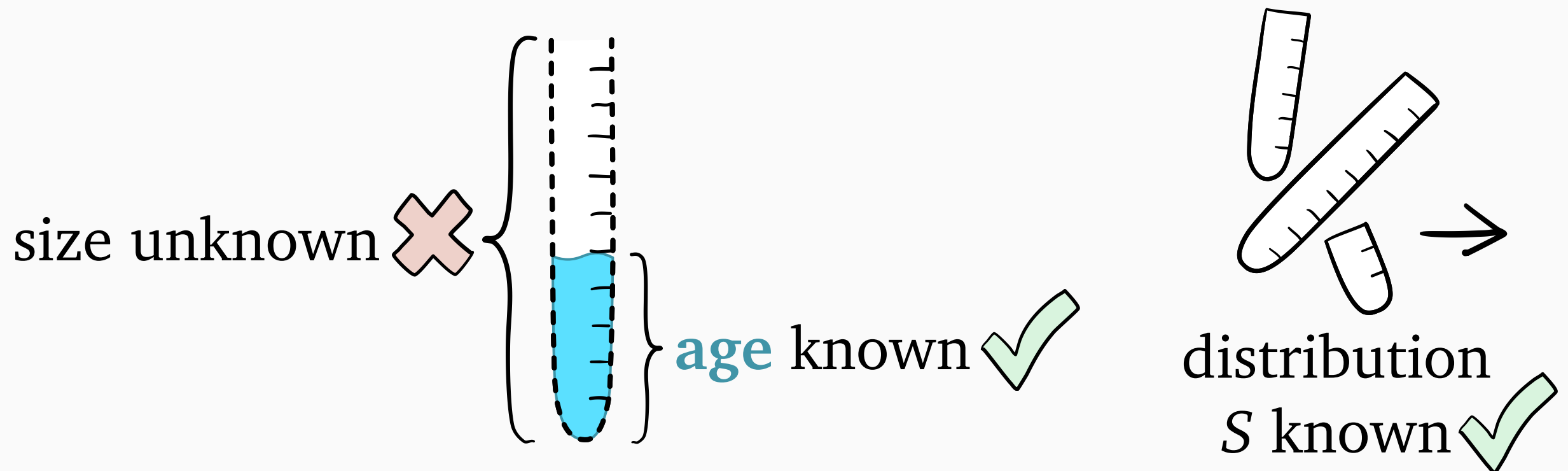


# “Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes

## SERPT

shortest *expected* remaining processing time



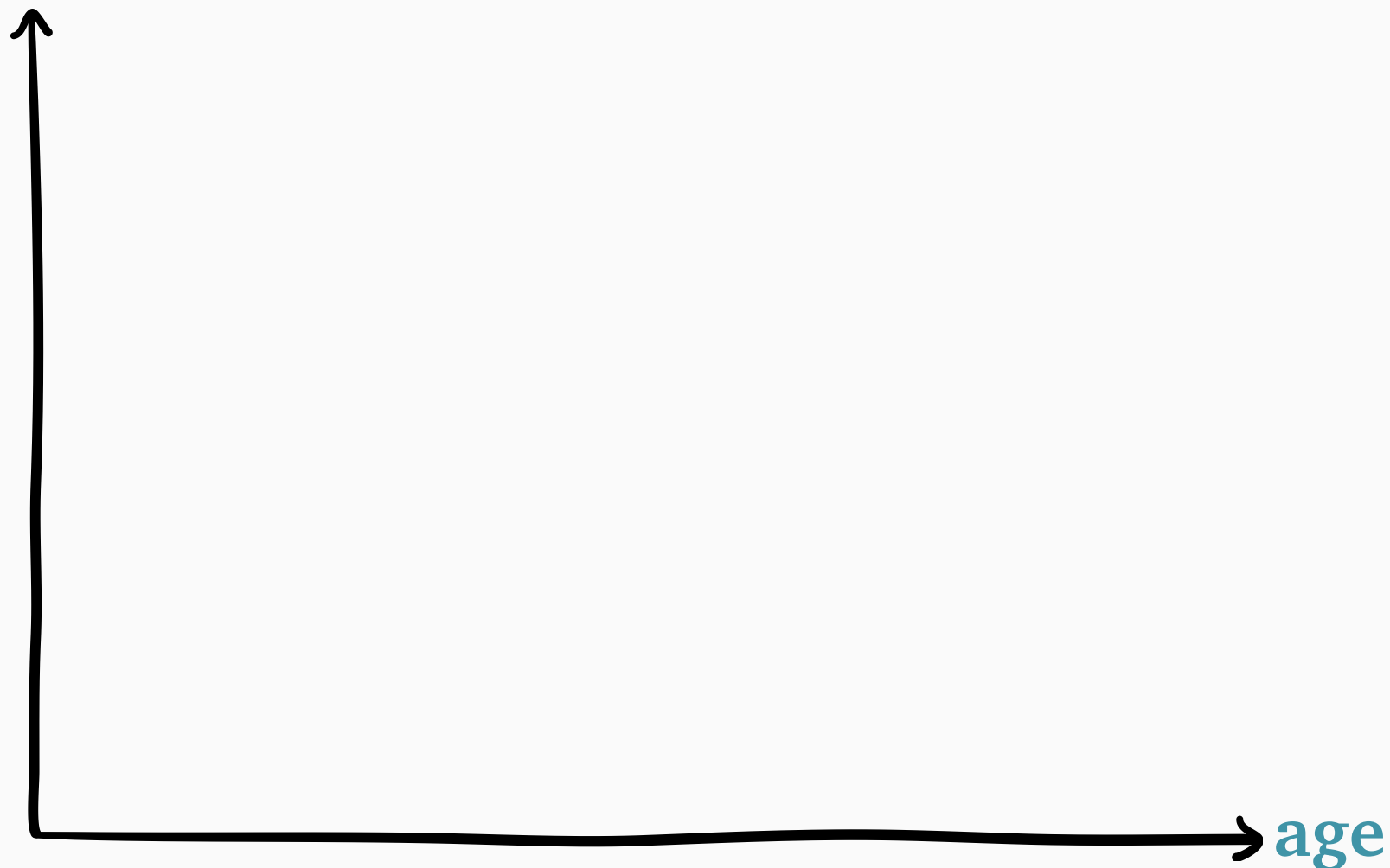
“Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes

## SERPT

shortest *expected* remaining processing time

rank



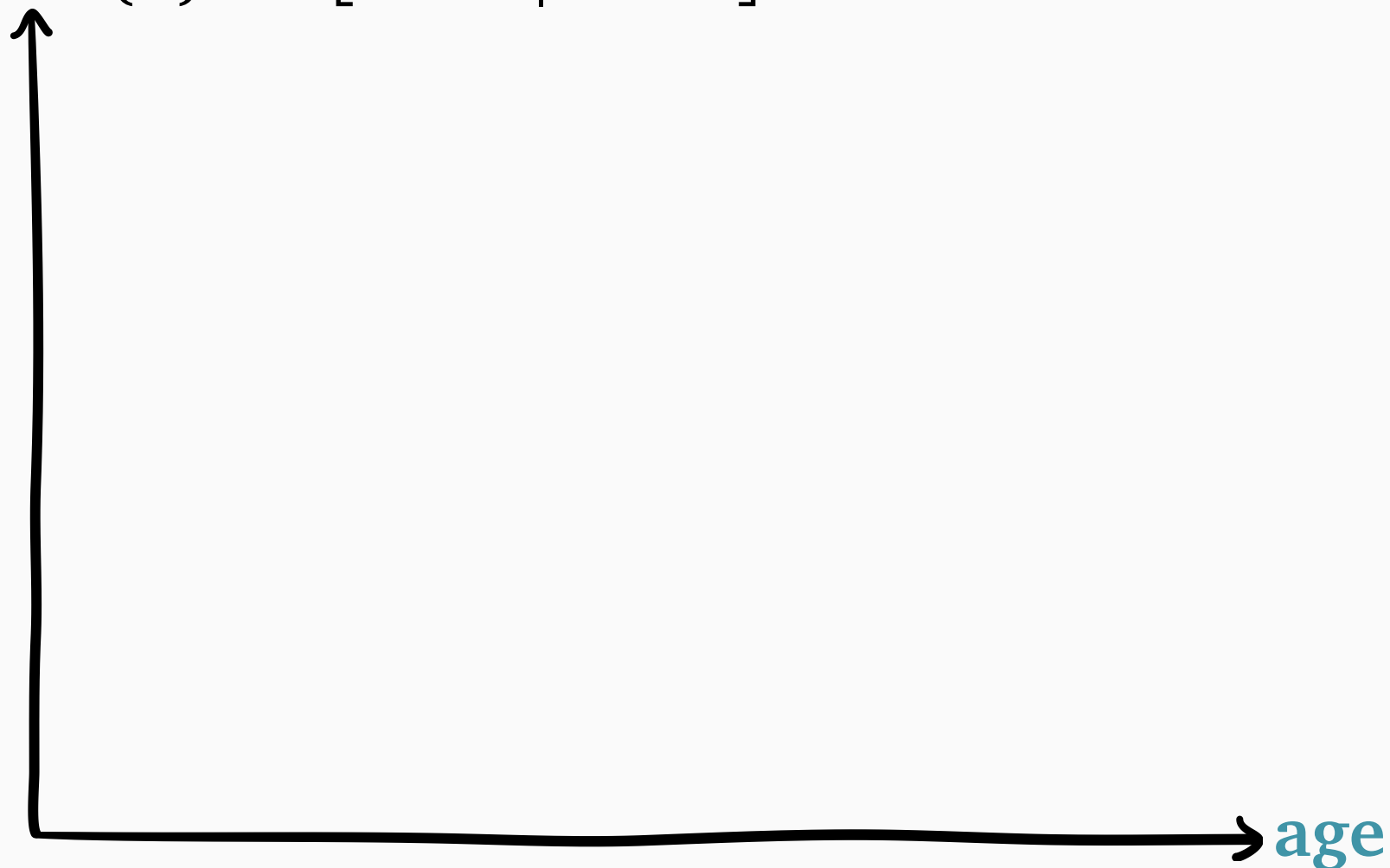
# “Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes

## SERPT

shortest *expected* remaining processing time

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



# “Complex” = *nonmonotonic rank*

Motivation:  
unknown job sizes

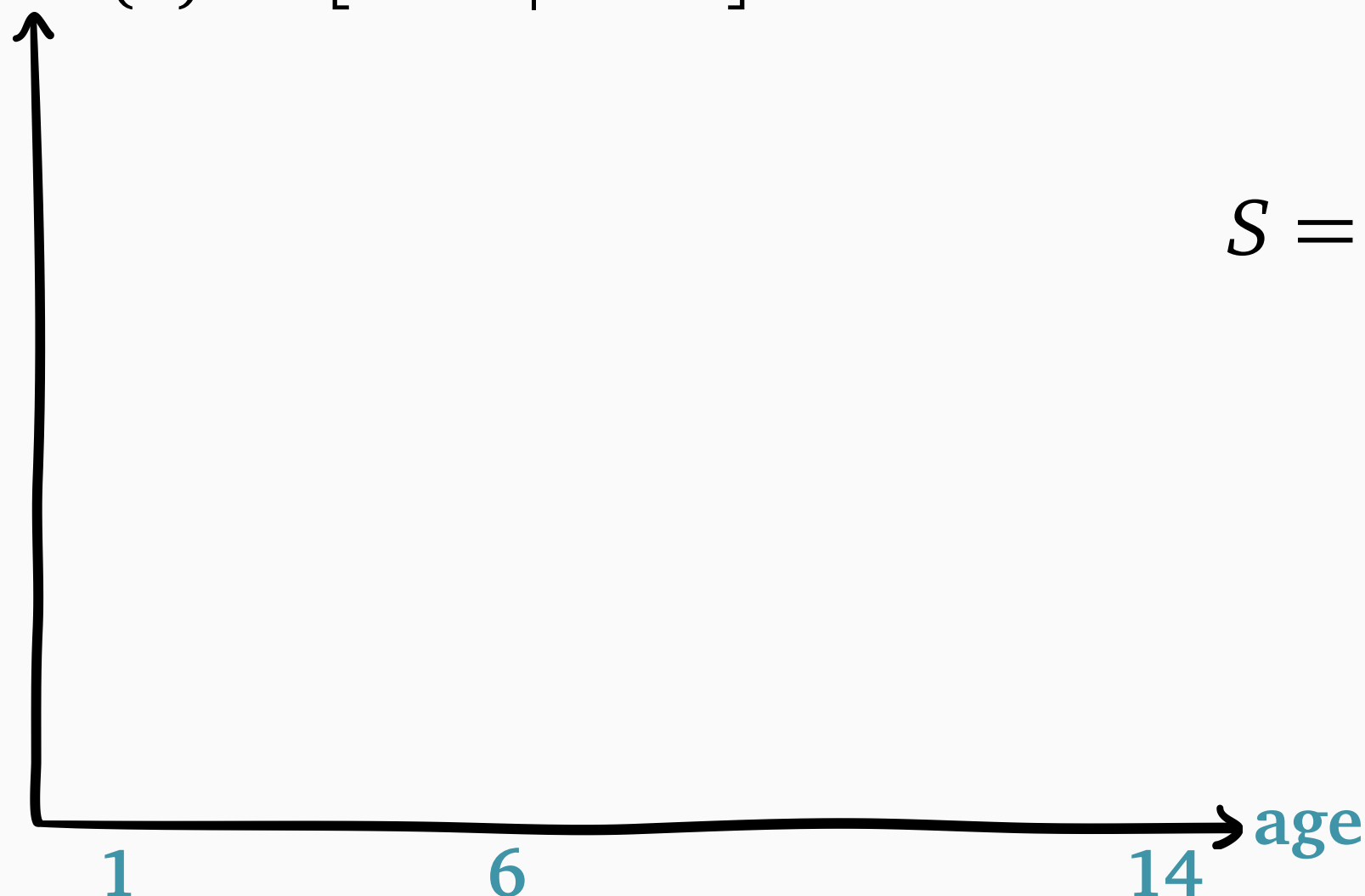
## SERPT

shortest *expected* remaining processing time

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

Job size distribution:

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



# “Complex” = *nonmonotonic* rank

Motivation:  
unknown job sizes

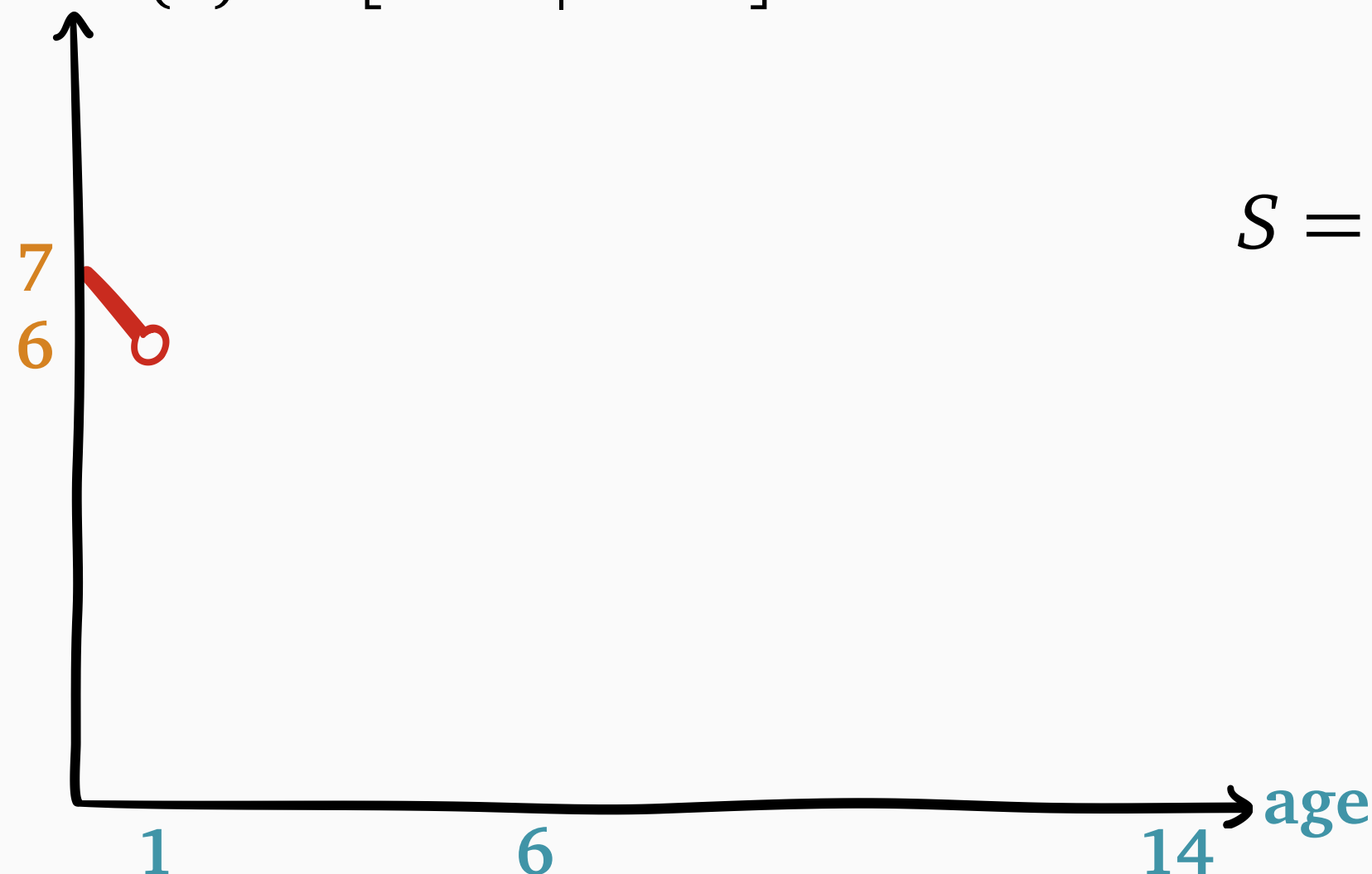
## SERPT

shortest *expected* remaining processing time

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

Job size distribution:

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



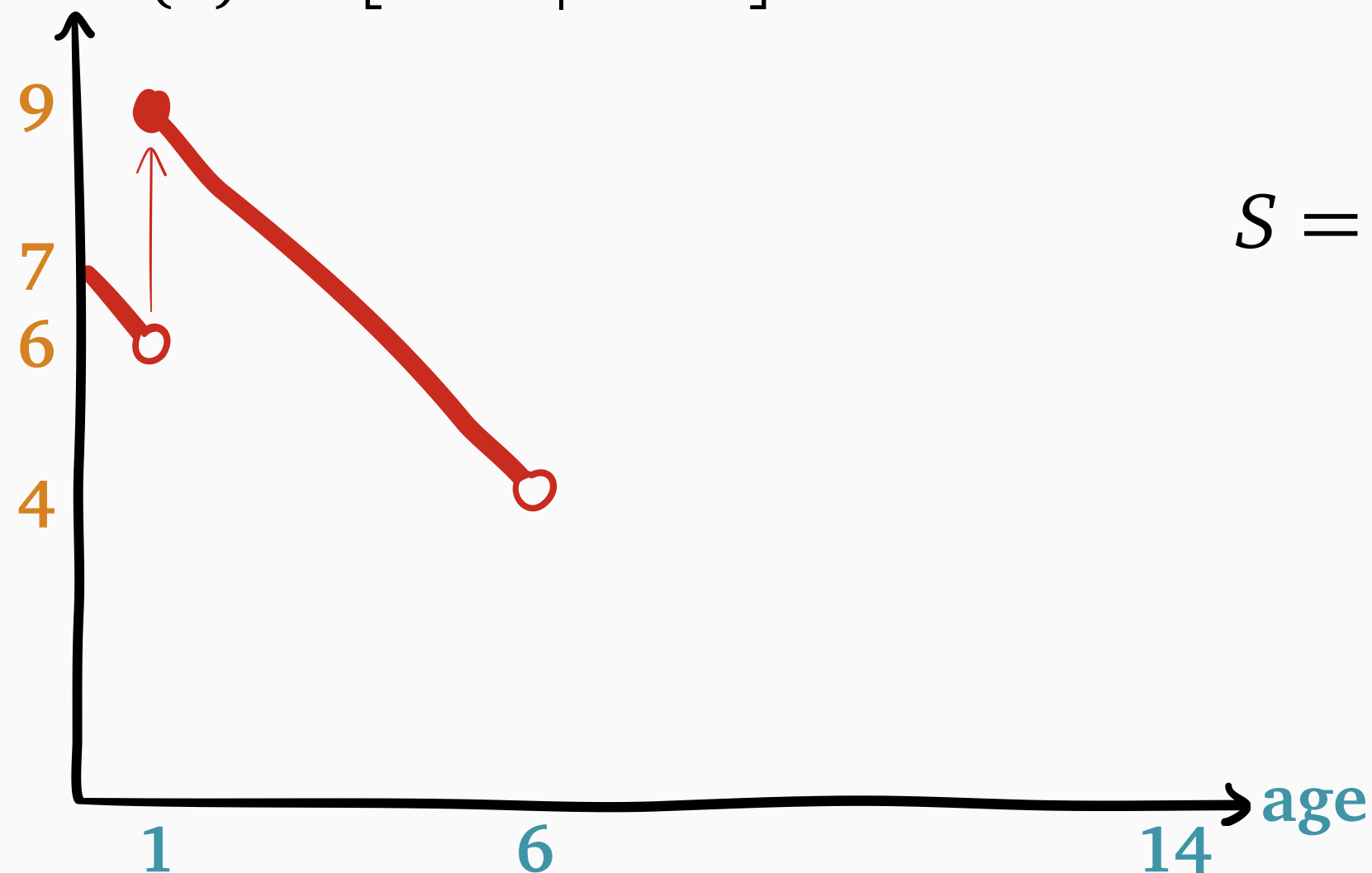
# “Complex” = *nonmonotonic rank*

Motivation:  
unknown job sizes

## SERPT

shortest *expected* remaining processing time

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



Job size distribution:

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

# “Complex” = *nonmonotonic rank*

Motivation:  
unknown job sizes

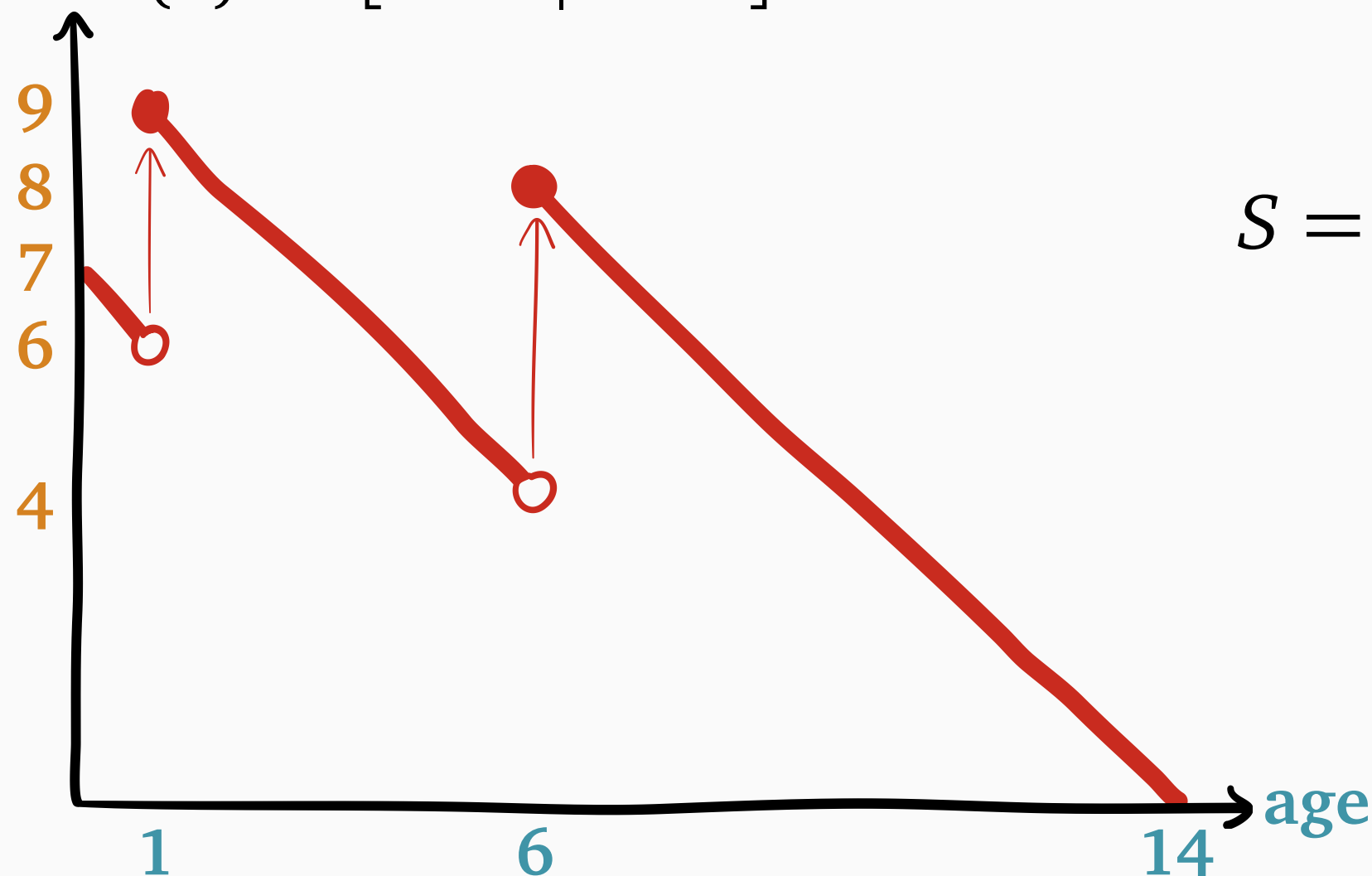
## SERPT

shortest *expected* remaining processing time

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

Job size distribution:

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



# “Complex” = *nonmonotonic rank*

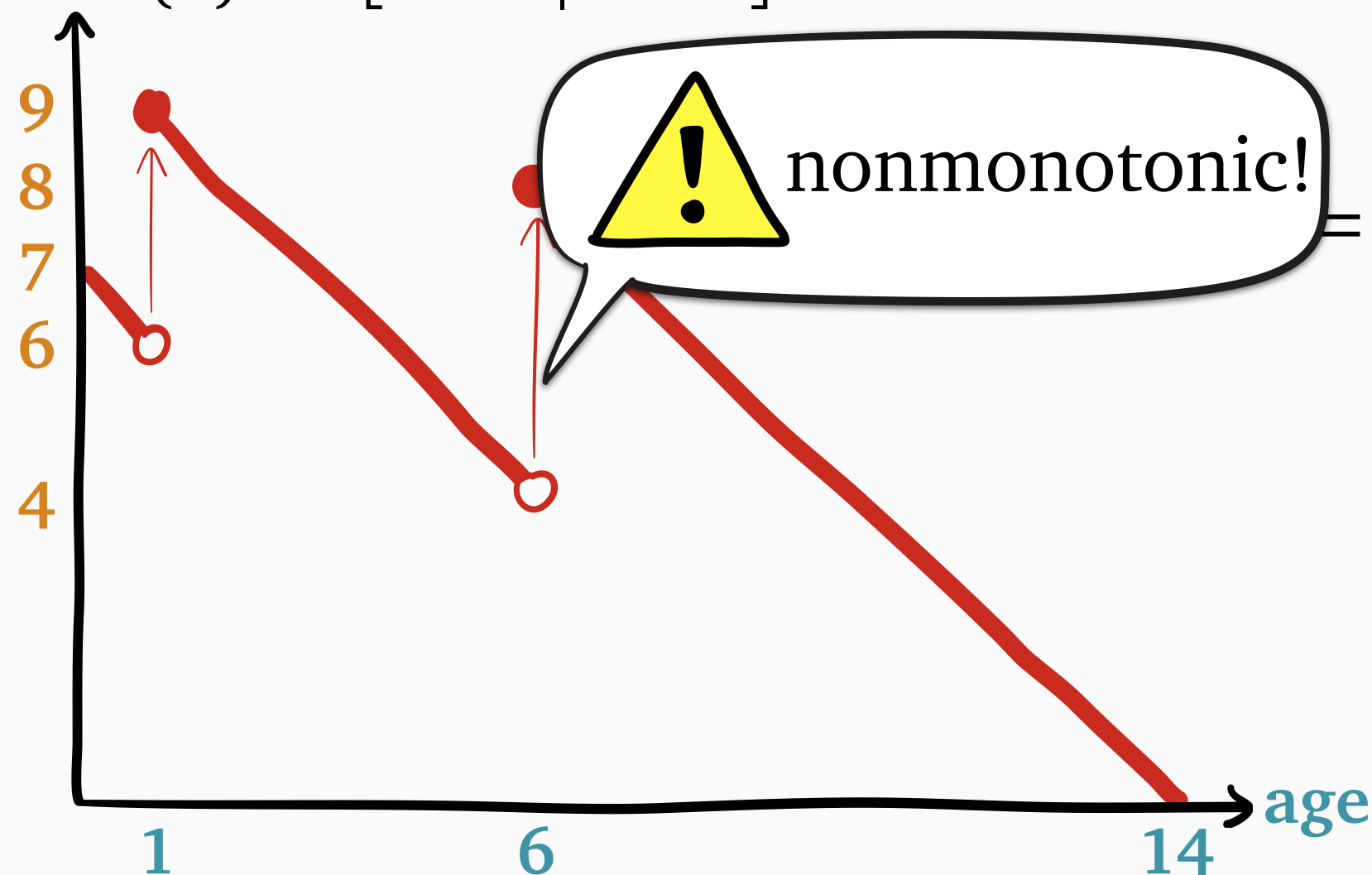
Motivation:  
unknown job sizes

## SERPT

shortest *expected* remaining processing time

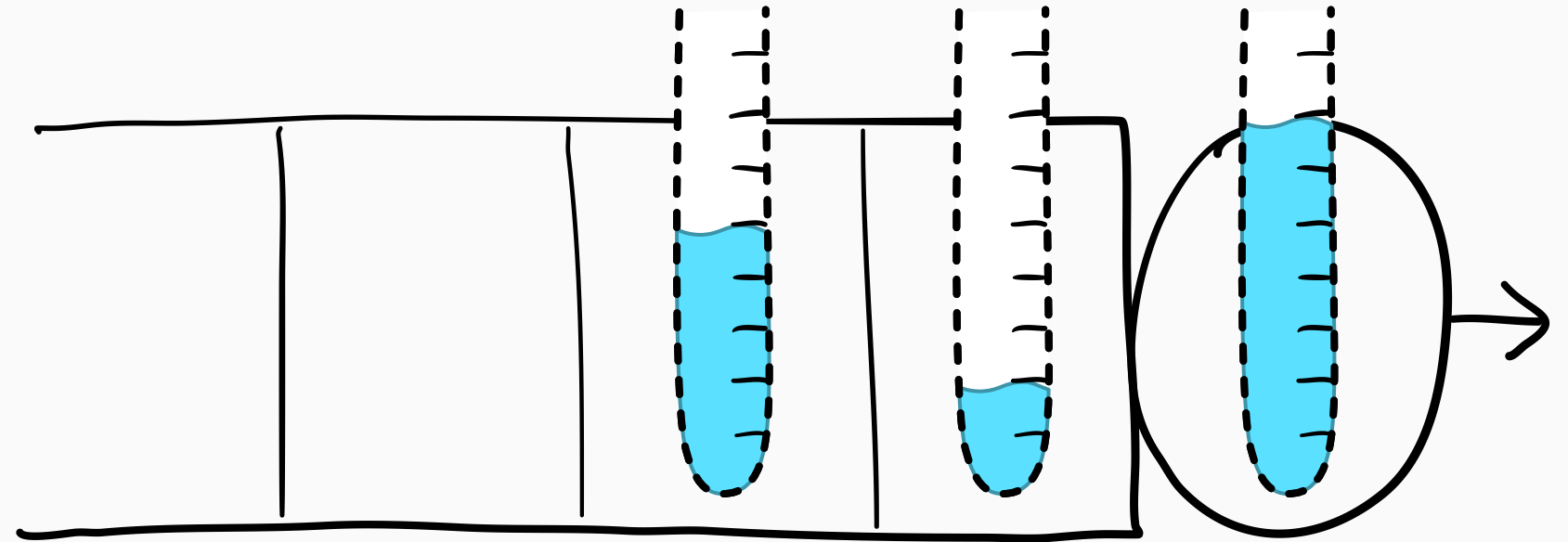
$$\text{rank}(a) = \mathbf{E}[S - a \mid S > a]$$

Job size distribution:

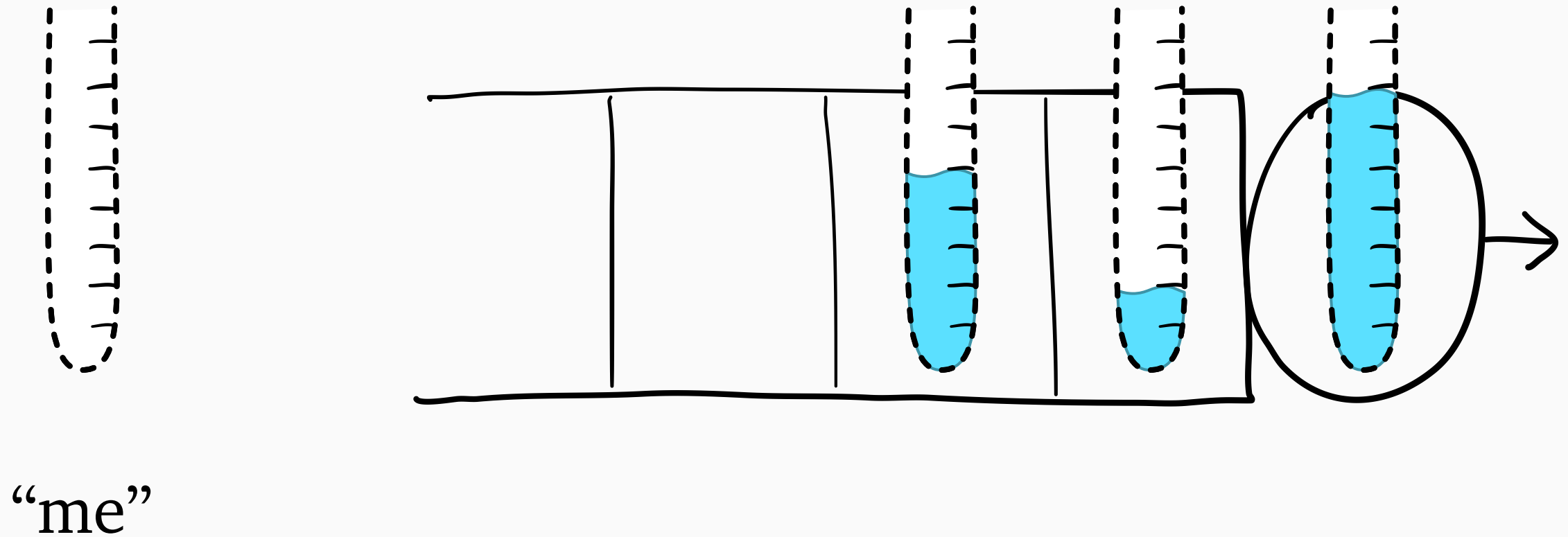




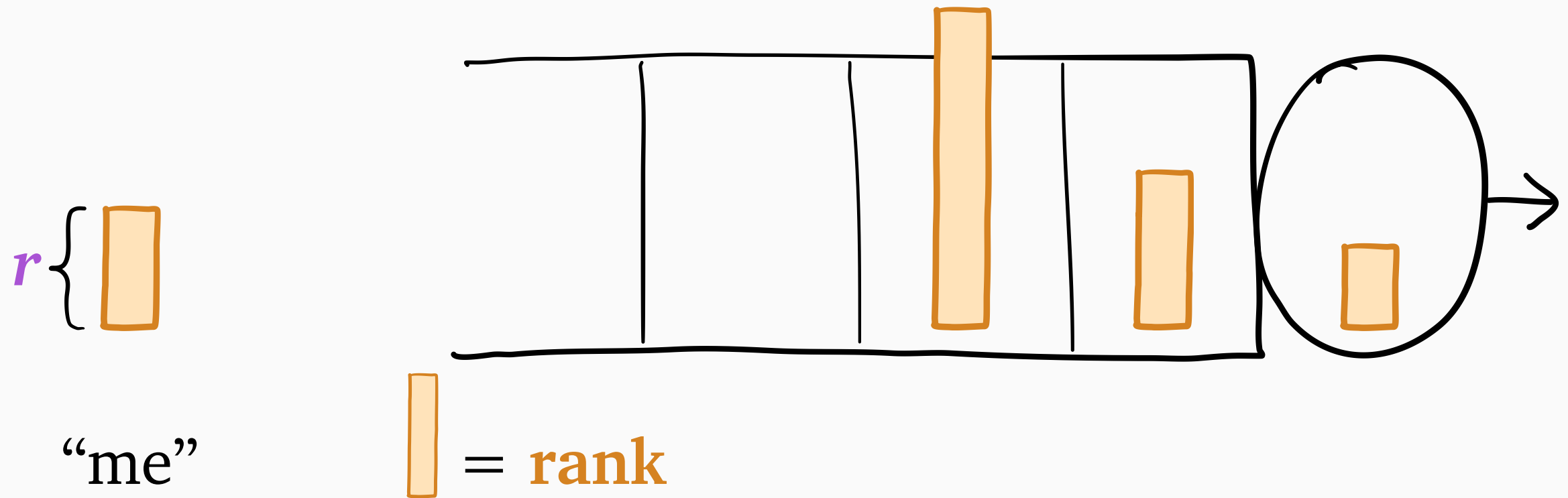
# Nonmonotonic **rank** is hard



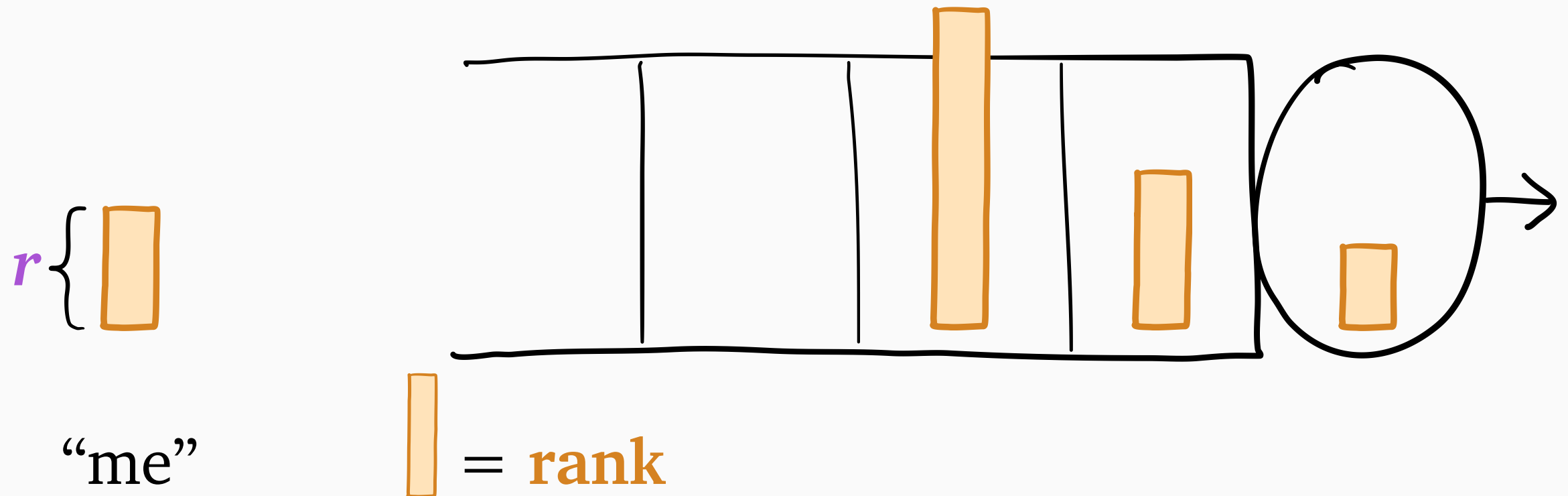
# Nonmonotonic **rank** is hard



# Nonmonotonic **rank** is hard



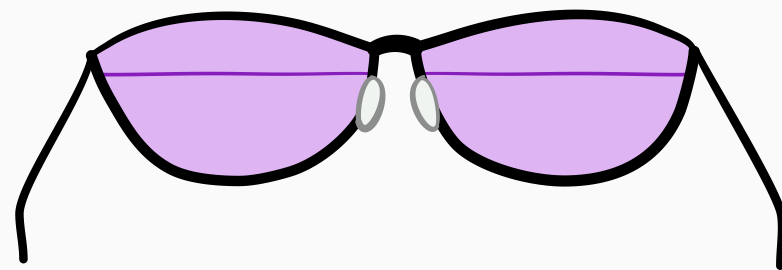
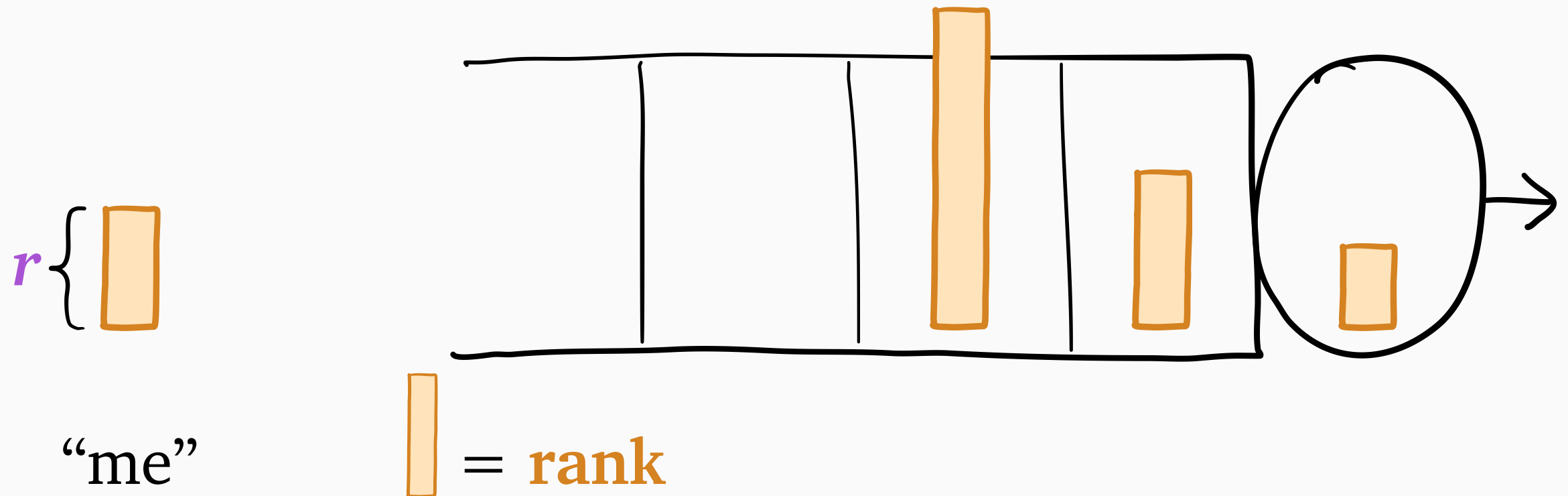
# Nonmonotonic **rank** is hard



Key quantity: observed  $r$ -work  $W(r)$

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

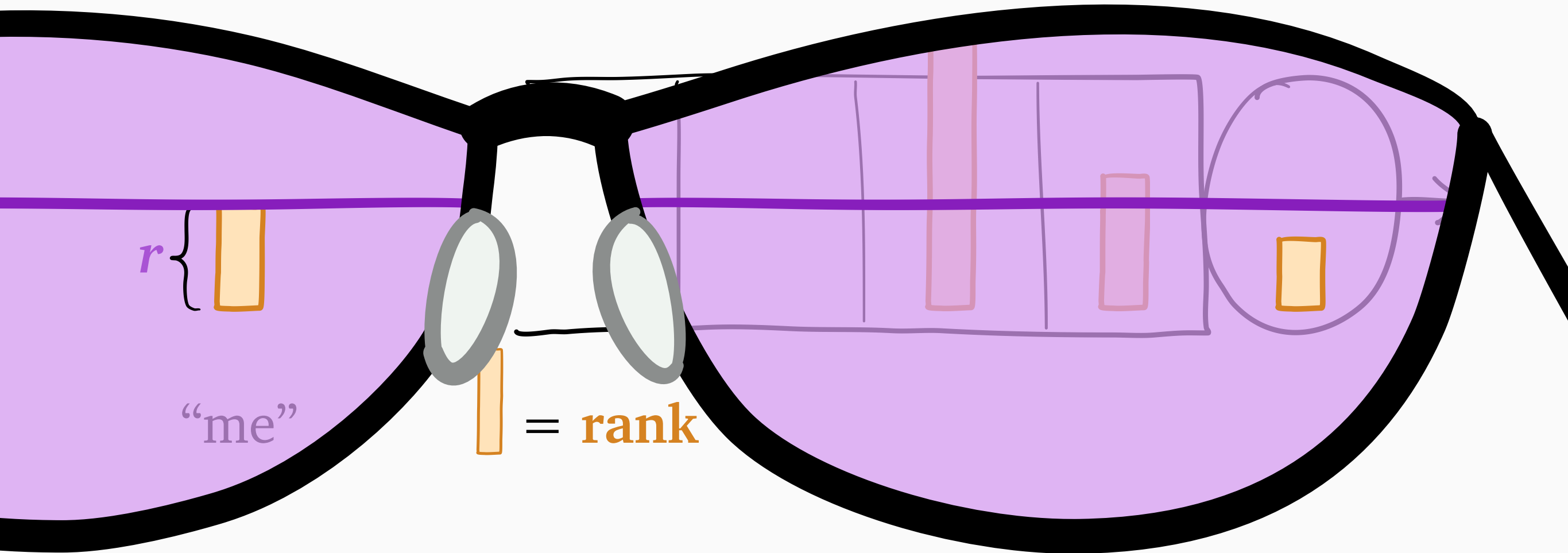
# Nonmonotonic **rank** is hard



Key quantity: observed  $r$ -work  $W(r)$

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

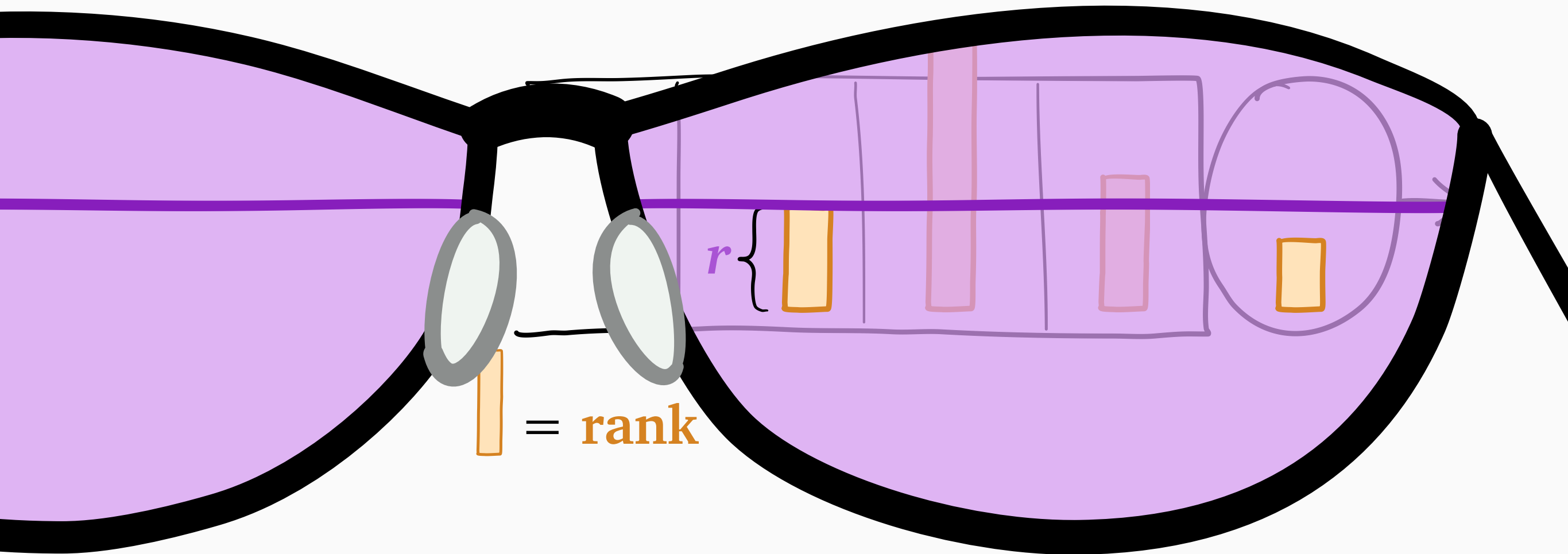
# Nonmonotonic **rank** is hard



Key quantity: observed  $r$ -work  $W(r)$

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

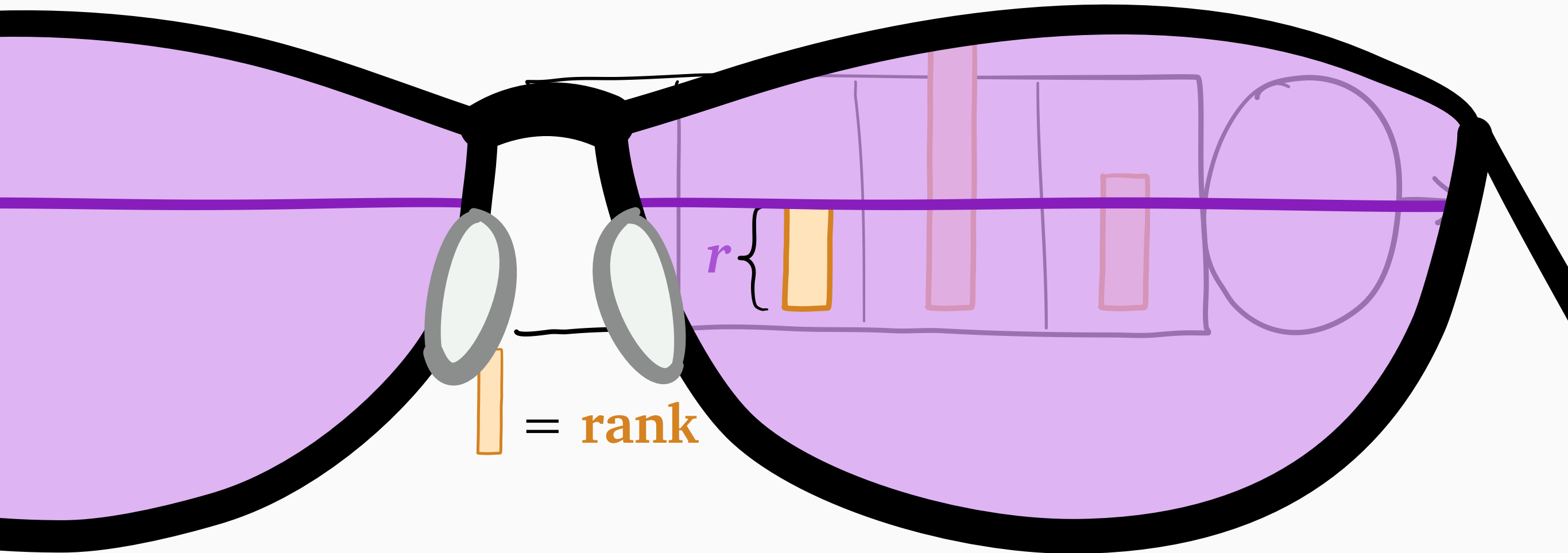
# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

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

# Nonmonotonic **rank** is hard

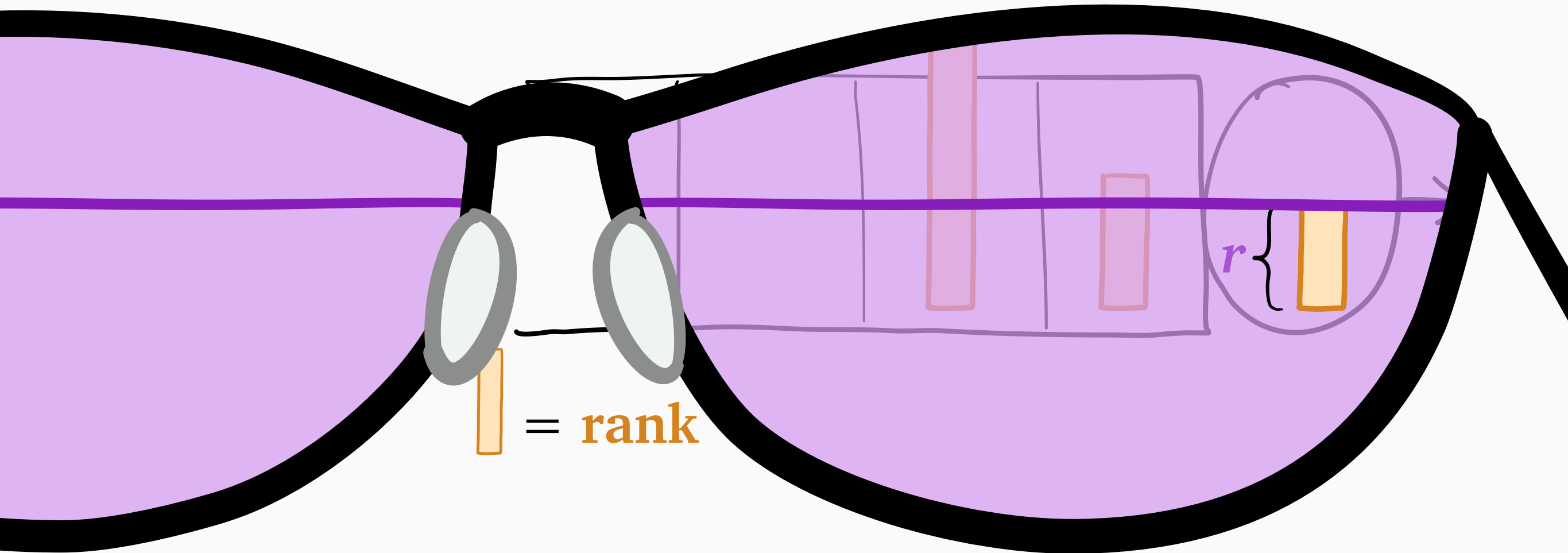


Key quantity: observed  $r$ -work  $W(r)$

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



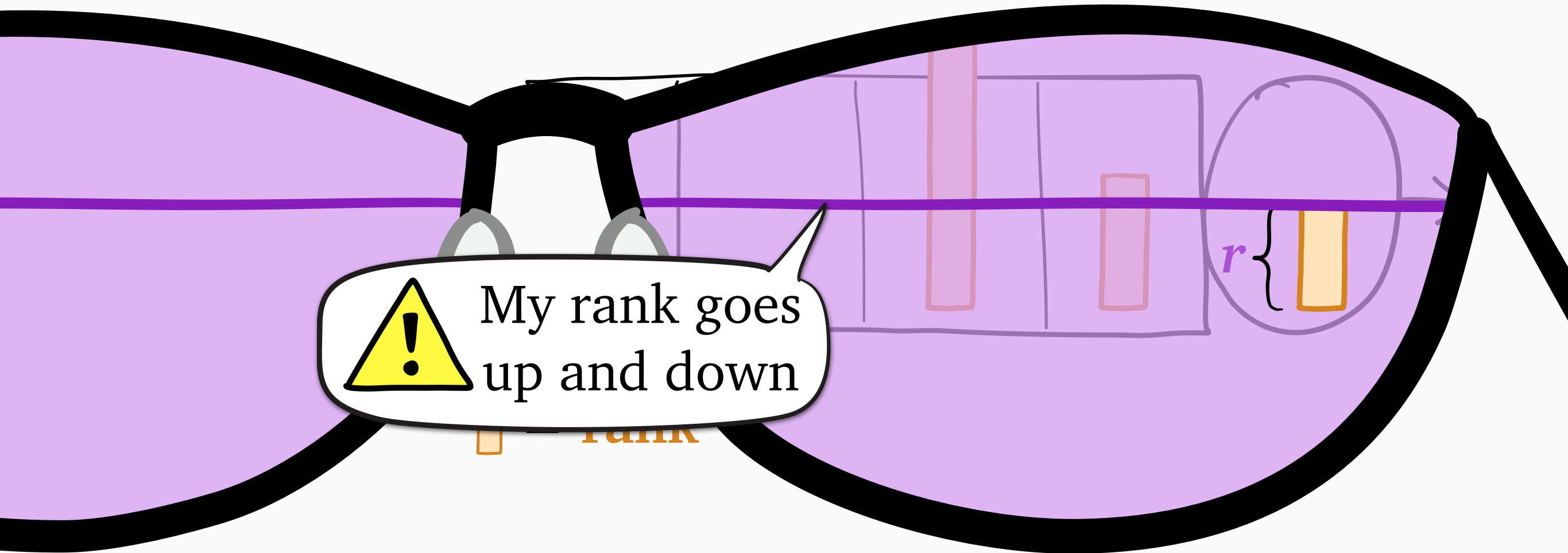
# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

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

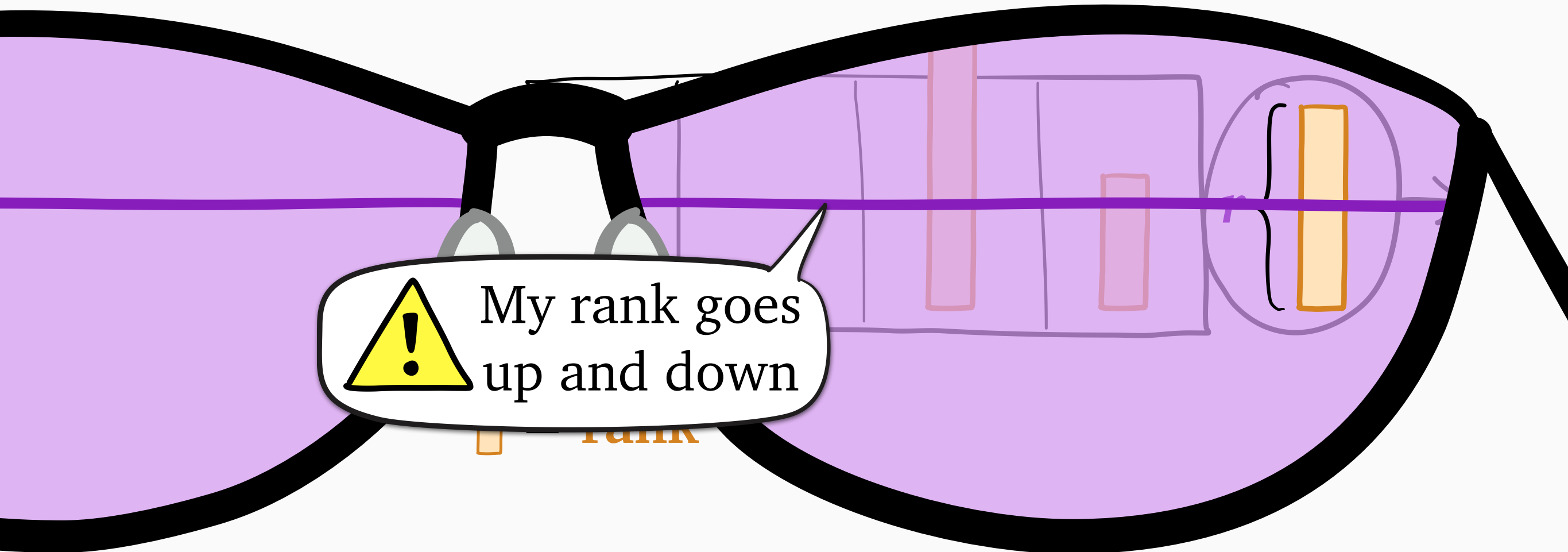
# Nonmonotonic **rank** is hard



Key quantity: observed  $r$ -work  $W(r)$

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

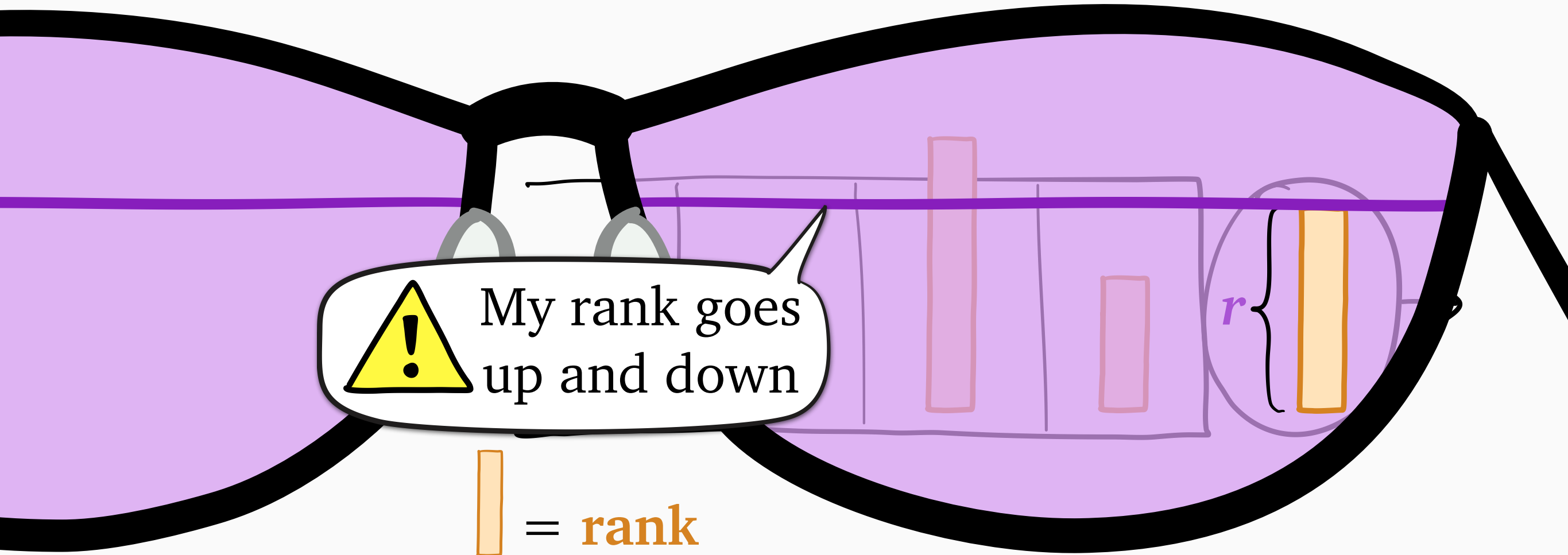
# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

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

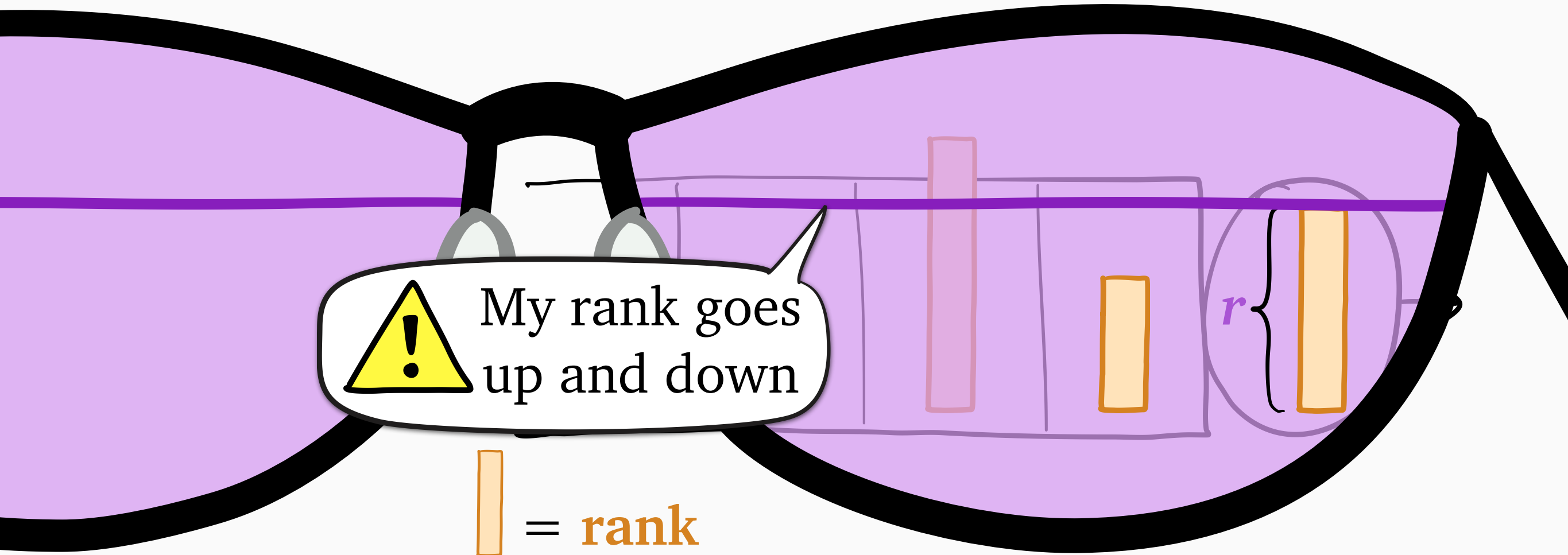
# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

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

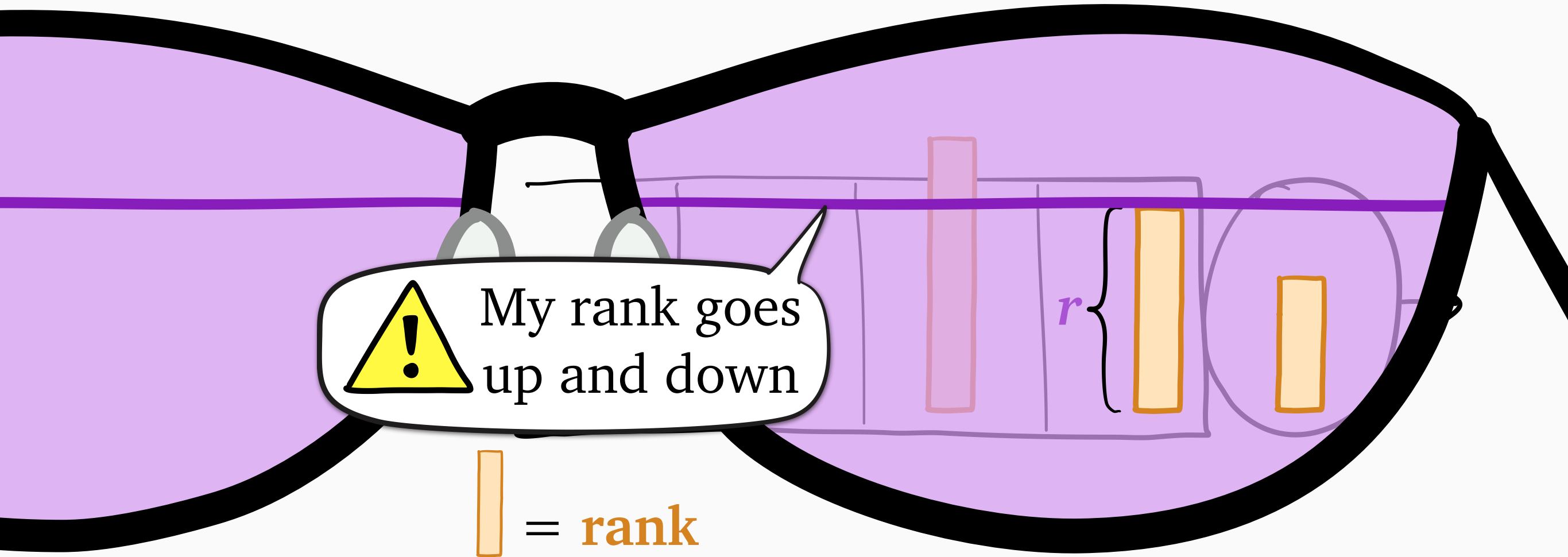
# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

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

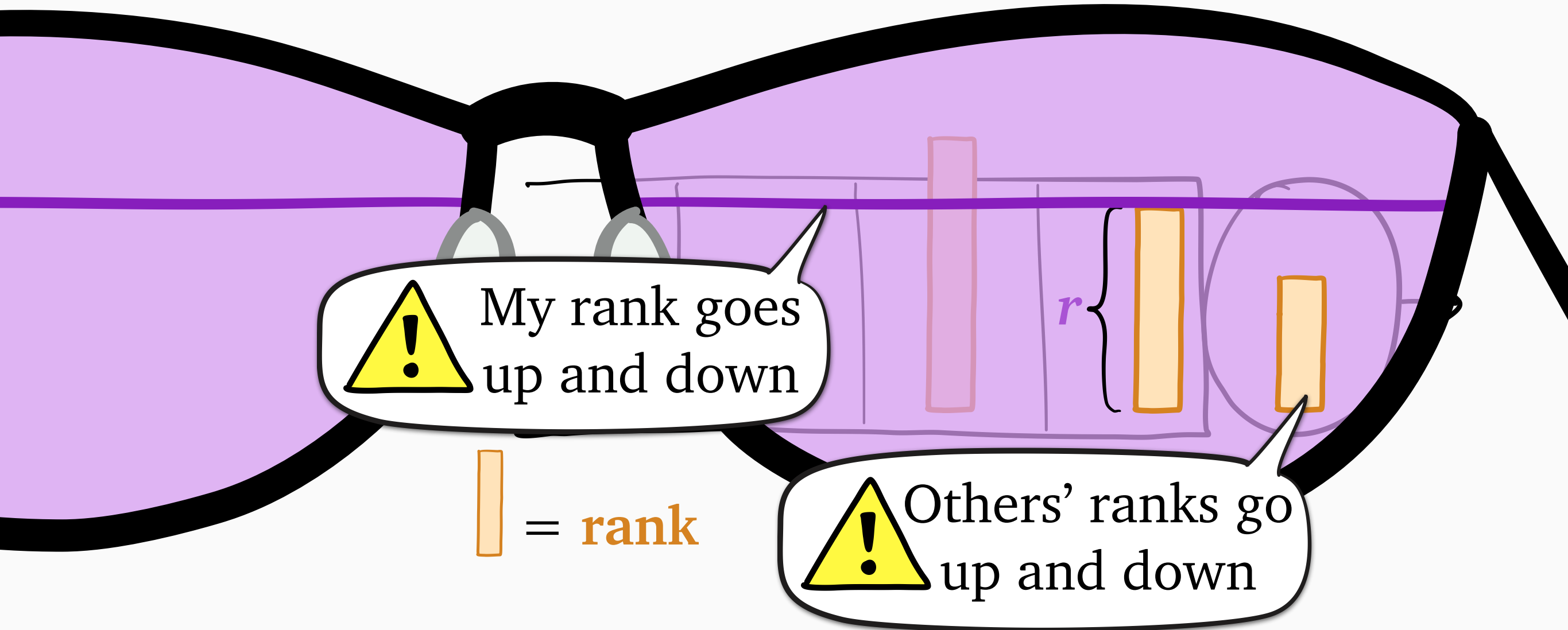
# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

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

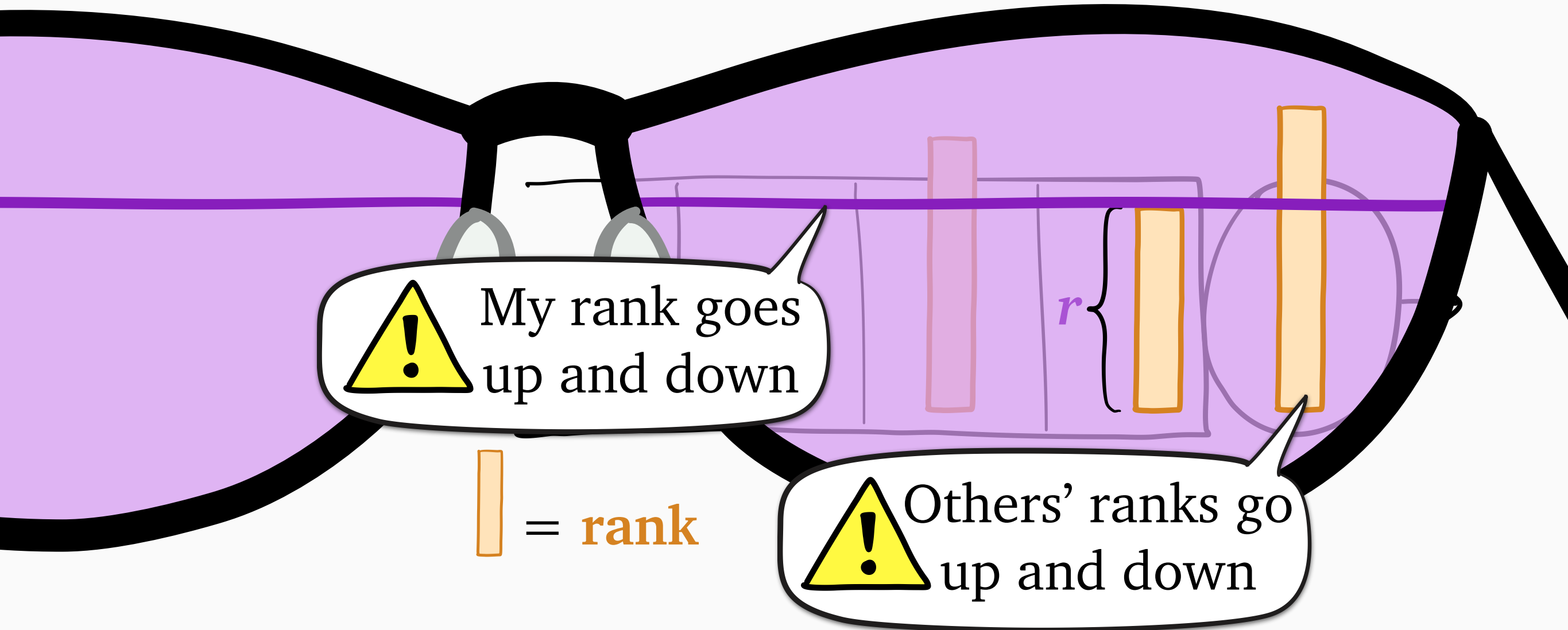
# Nonmonotonic **rank** is hard



Key quantity: observed  $r$ -work  $W(r)$

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

# Nonmonotonic **rank** is hard

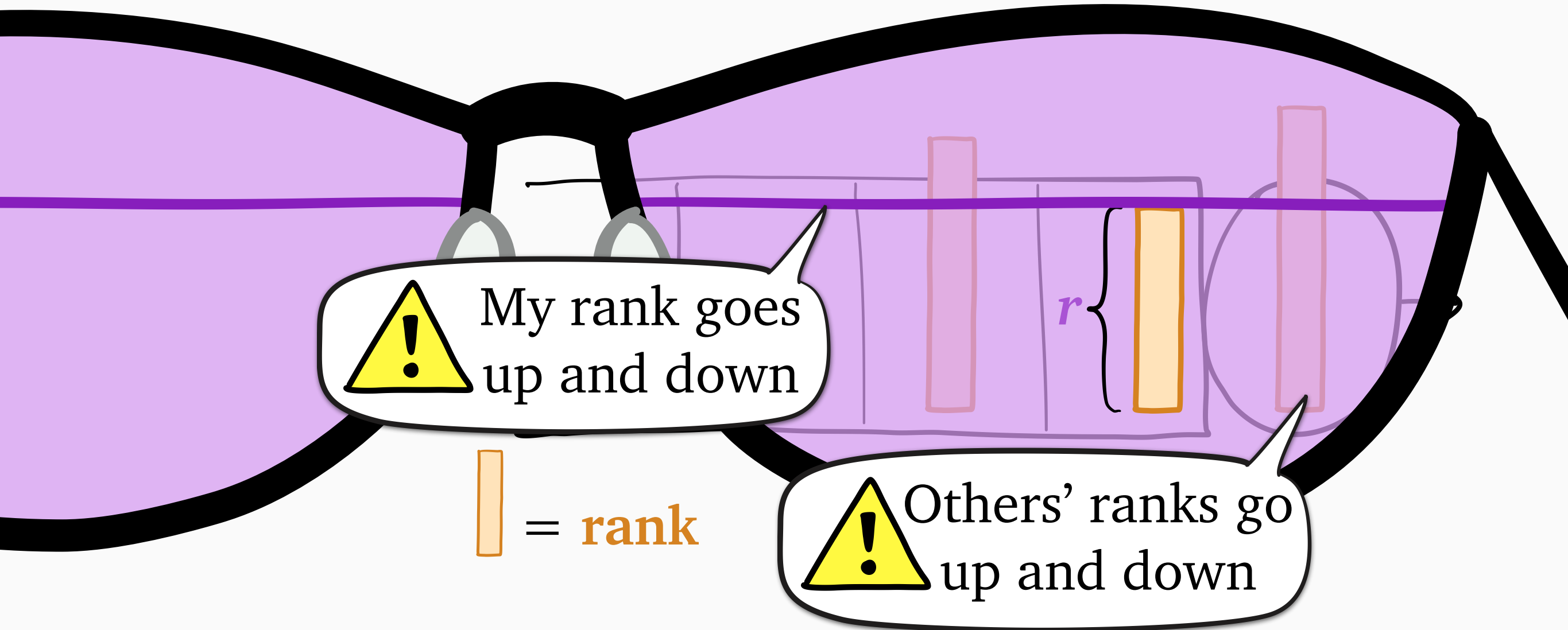


Key quantity: observed  $r$ -work  $W(r)$

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



# Nonmonotonic **rank** is hard



Key quantity: observed *r*-work  $W(r)$

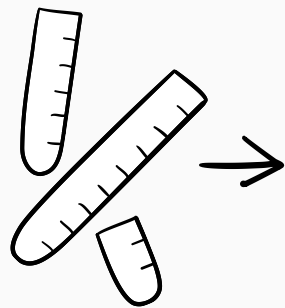
$W(r)$  = work relevant to job of **rank** *r*



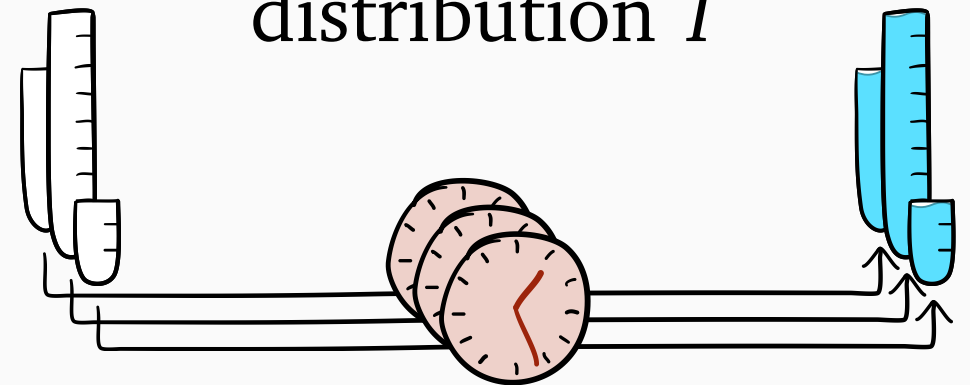
# SOAP

Schedule **O**rdered by **A**ge-based **P**riority

stochastic arrival  
process  $\lambda, S$



response time  
distribution  $T$

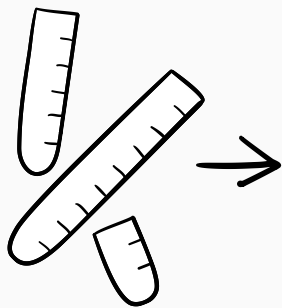




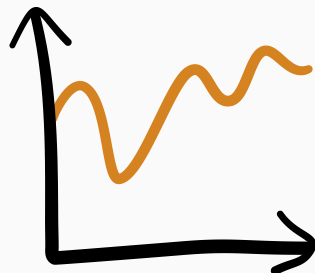
# SOAP

Schedule **O**rdered by **A**ge-based **P**riority

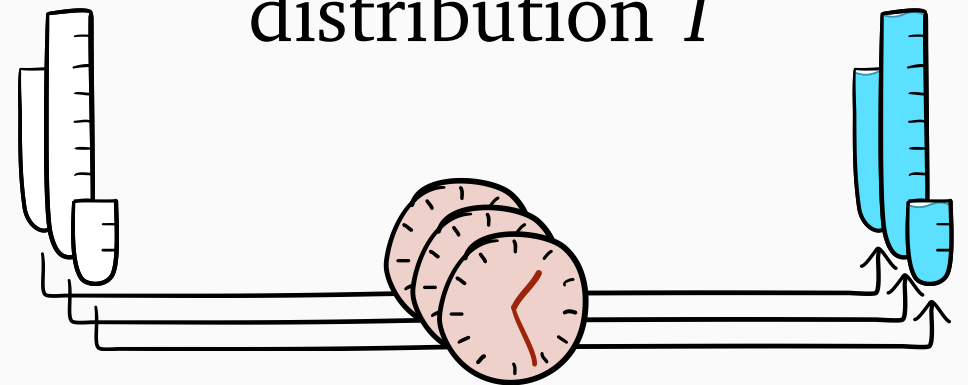
stochastic arrival  
process  $\lambda, S$



**rank** function



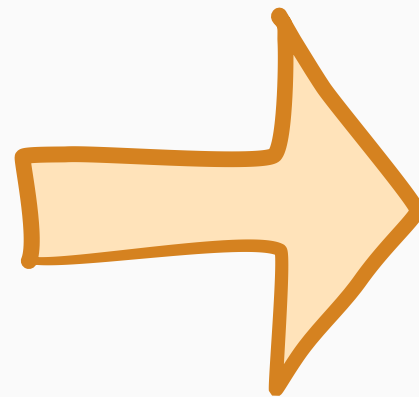
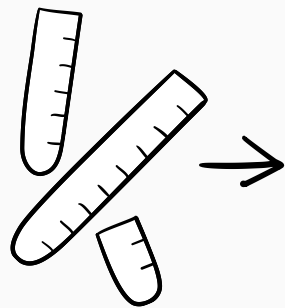
response time  
distribution  $T$



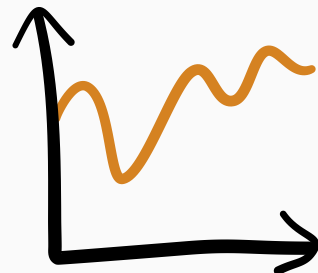
# SOAP

Schedule **O**rdered by **A**ge-based **P**riority

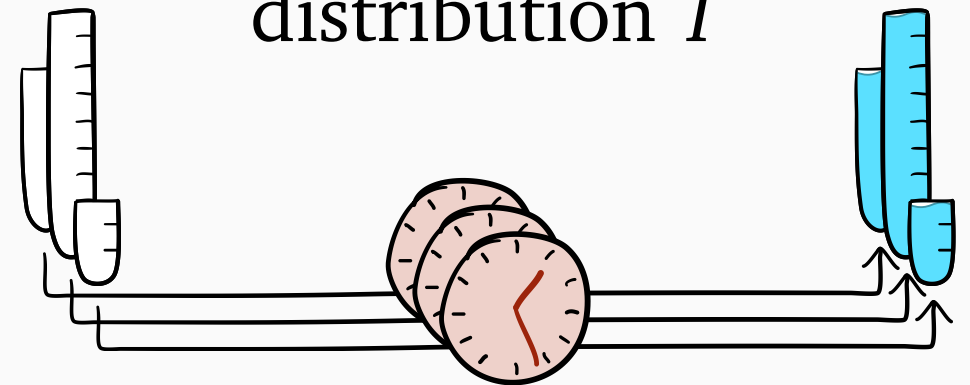
stochastic arrival  
process  $\lambda, S$



**rank** function



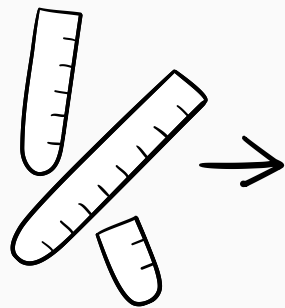
response time  
distribution  $T$



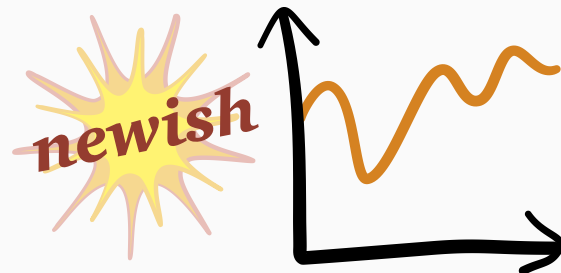
# SOAP

Schedule **O**rdered by **A**ge-based **P**riority

stochastic arrival  
process  $\lambda, S$

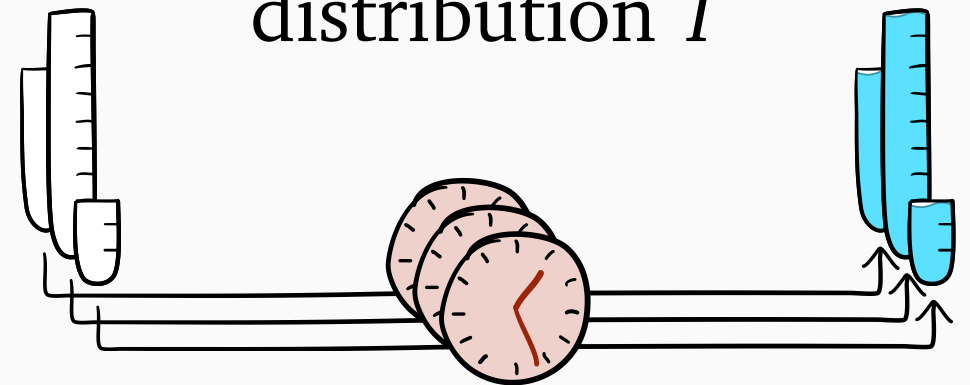


**rank** function

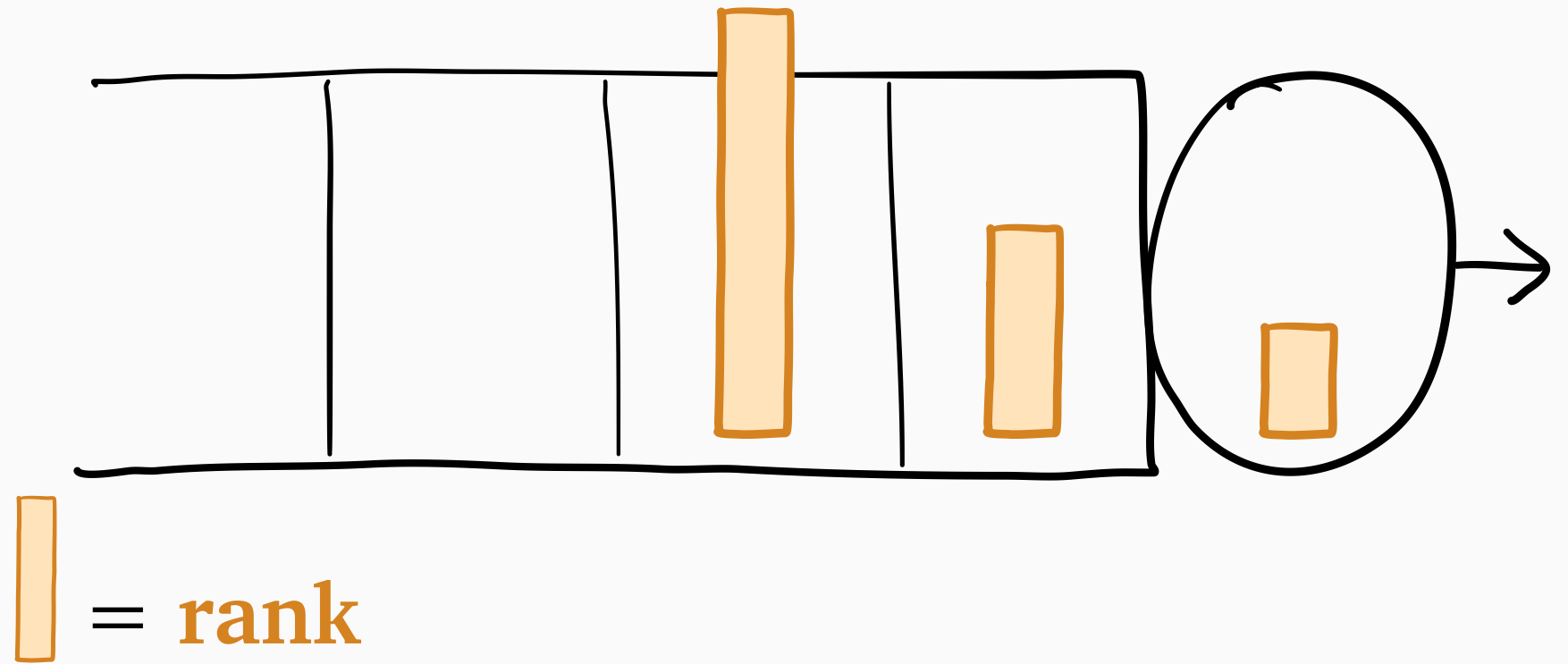


**NEW!**

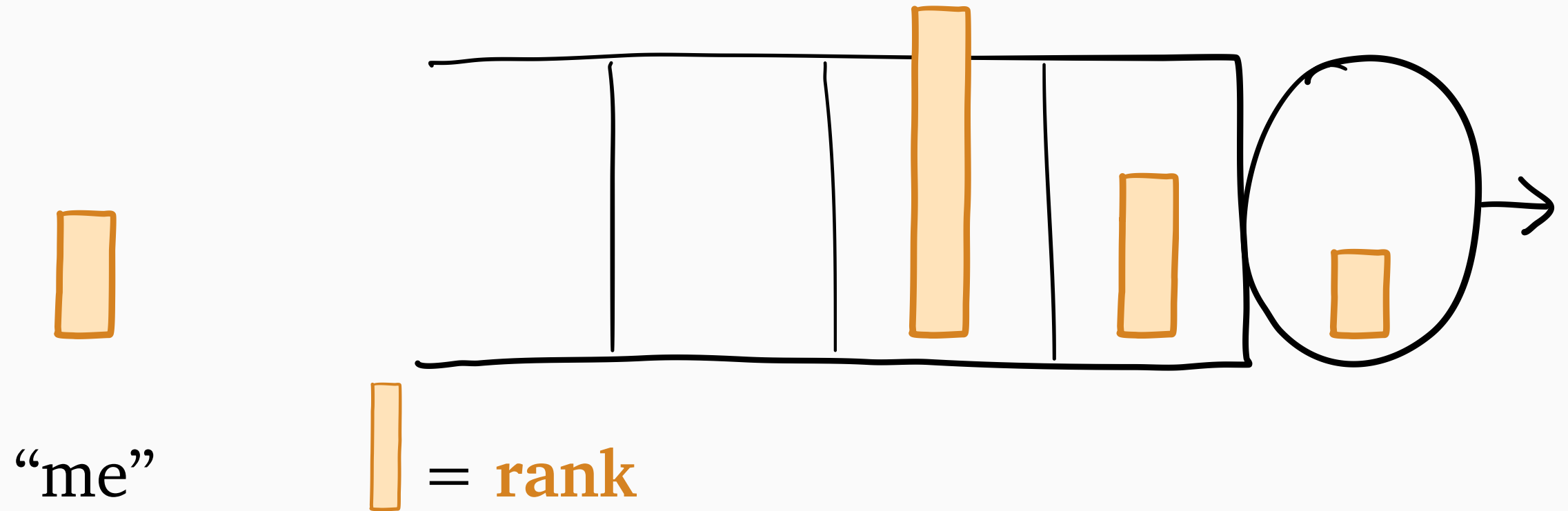
response time  
distribution  $T$



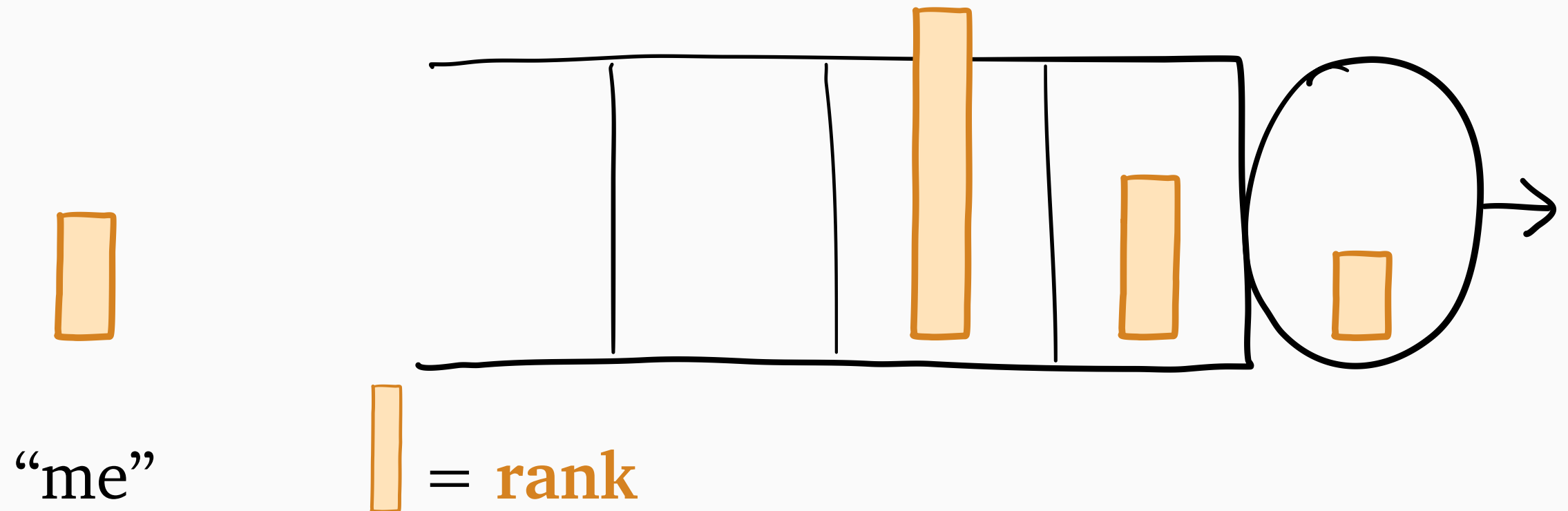
# Key SOAP insight



# Key SOAP insight



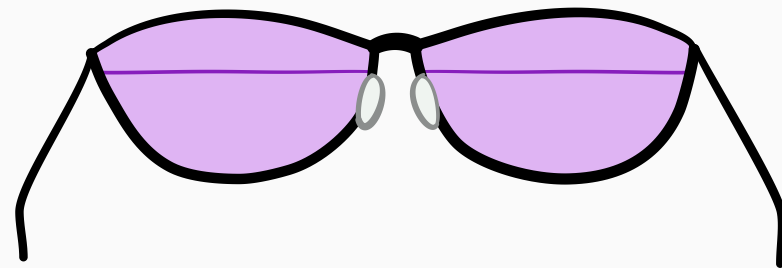
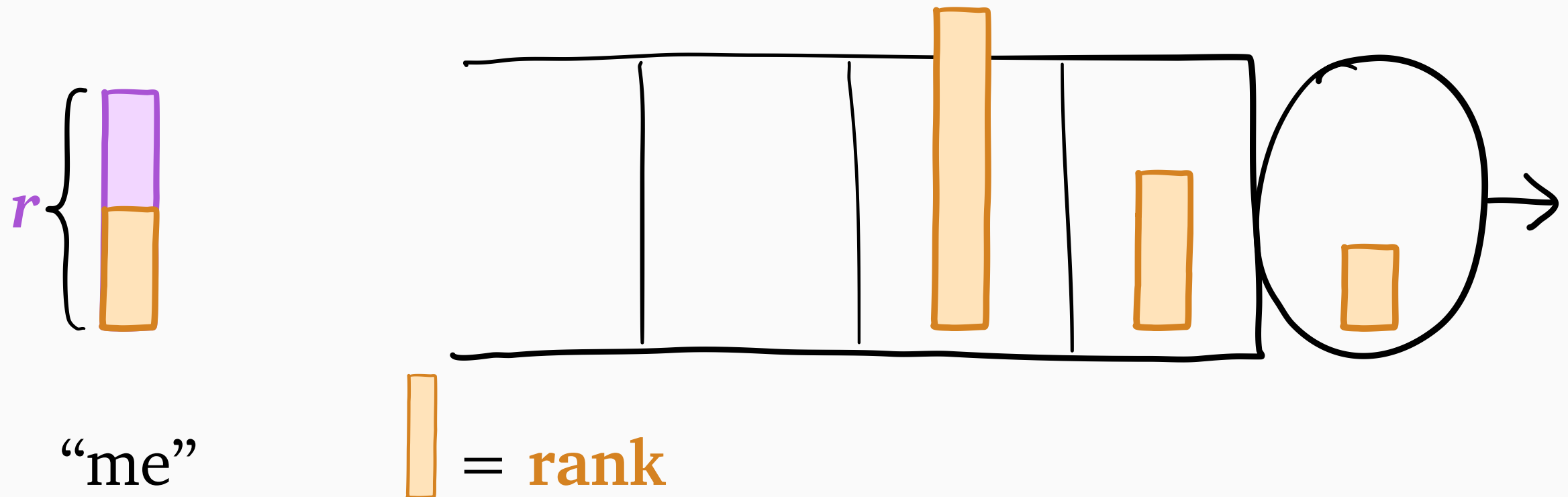
# Key SOAP insight



**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

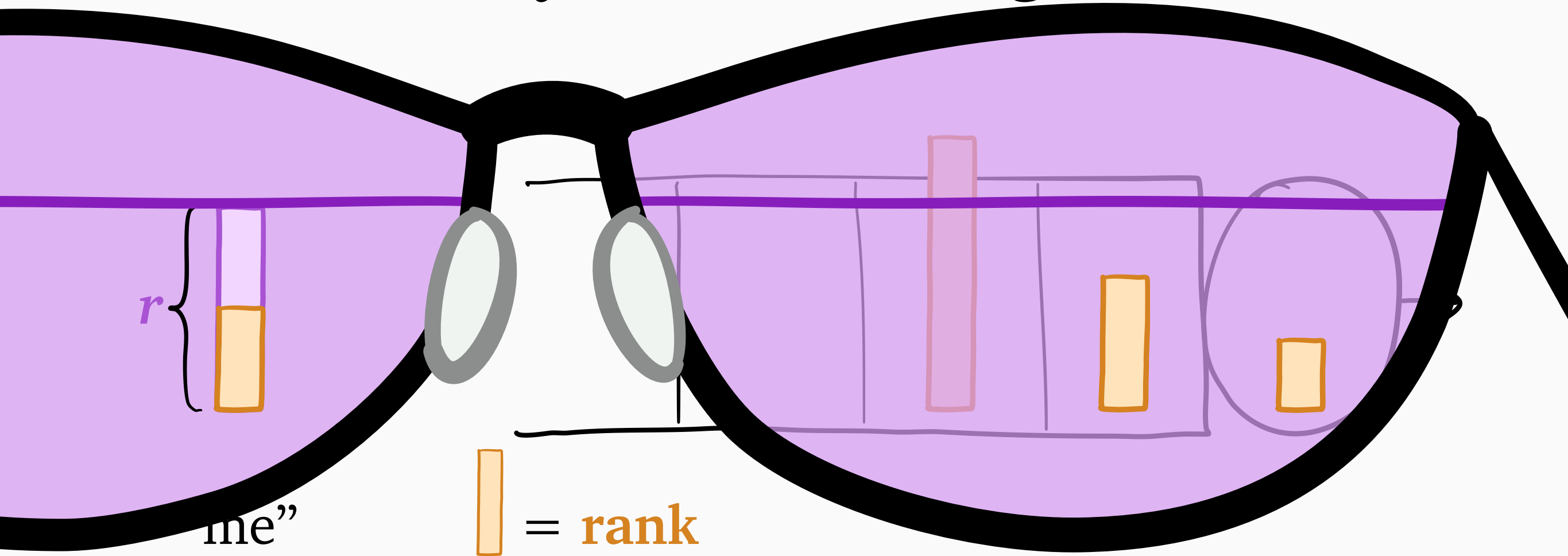


# Key SOAP insight



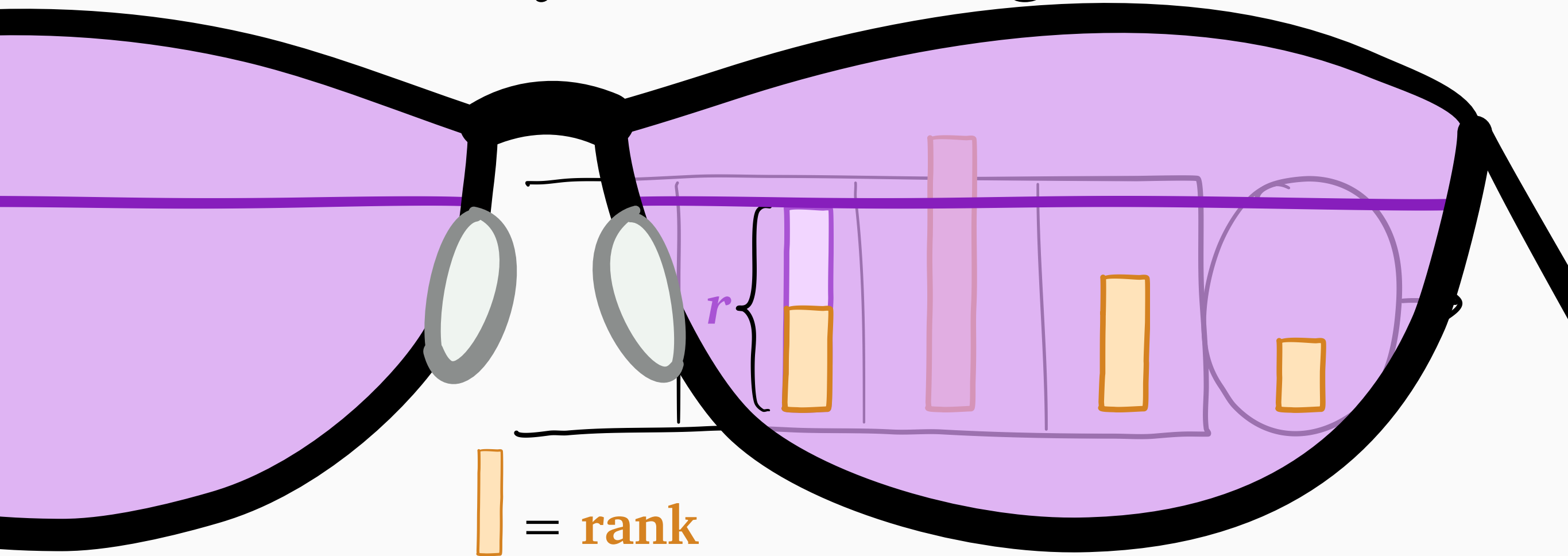
**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

# Key SOAP insight



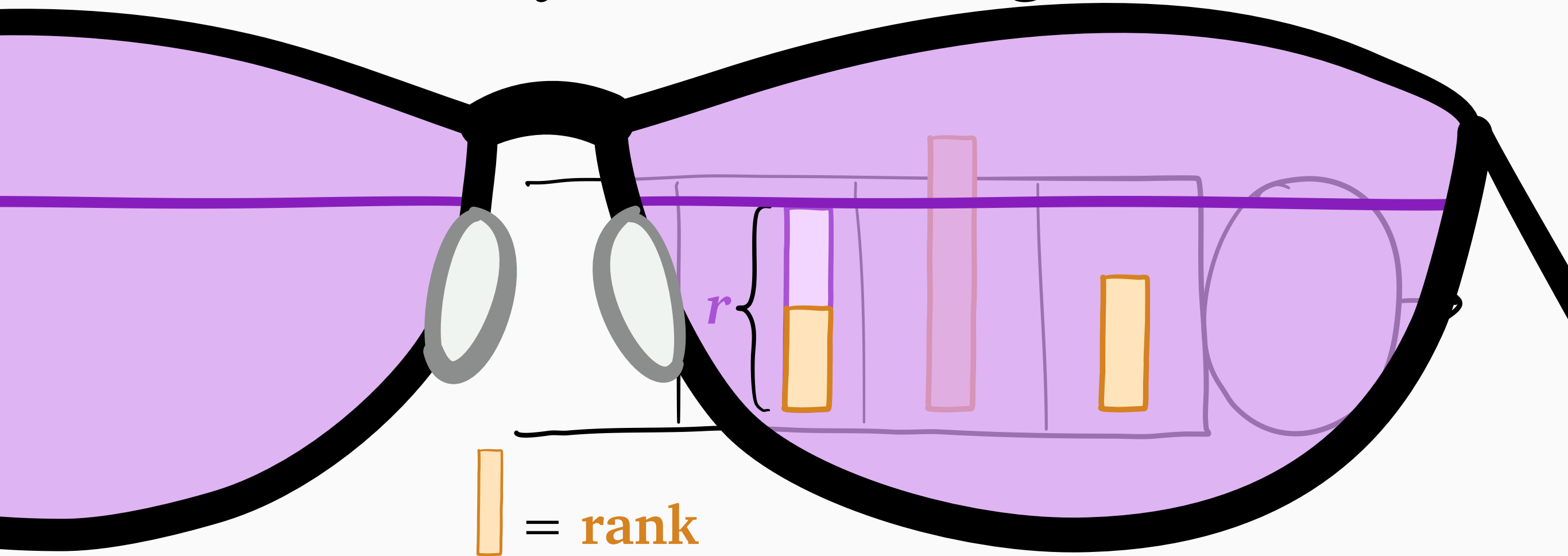
**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

# Key **SOAP** insight



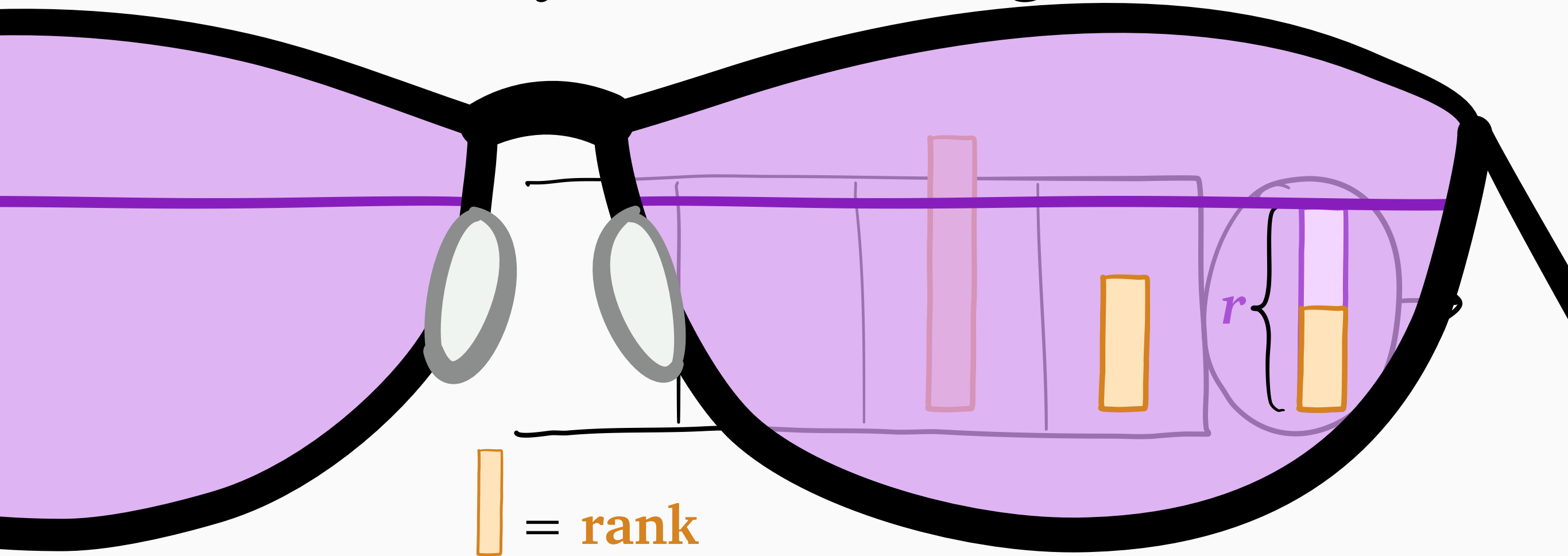
**Pessimism Principle:** compute  $W(r)$   
with  $r =$  my *worst future rank*

# Key **SOAP** insight



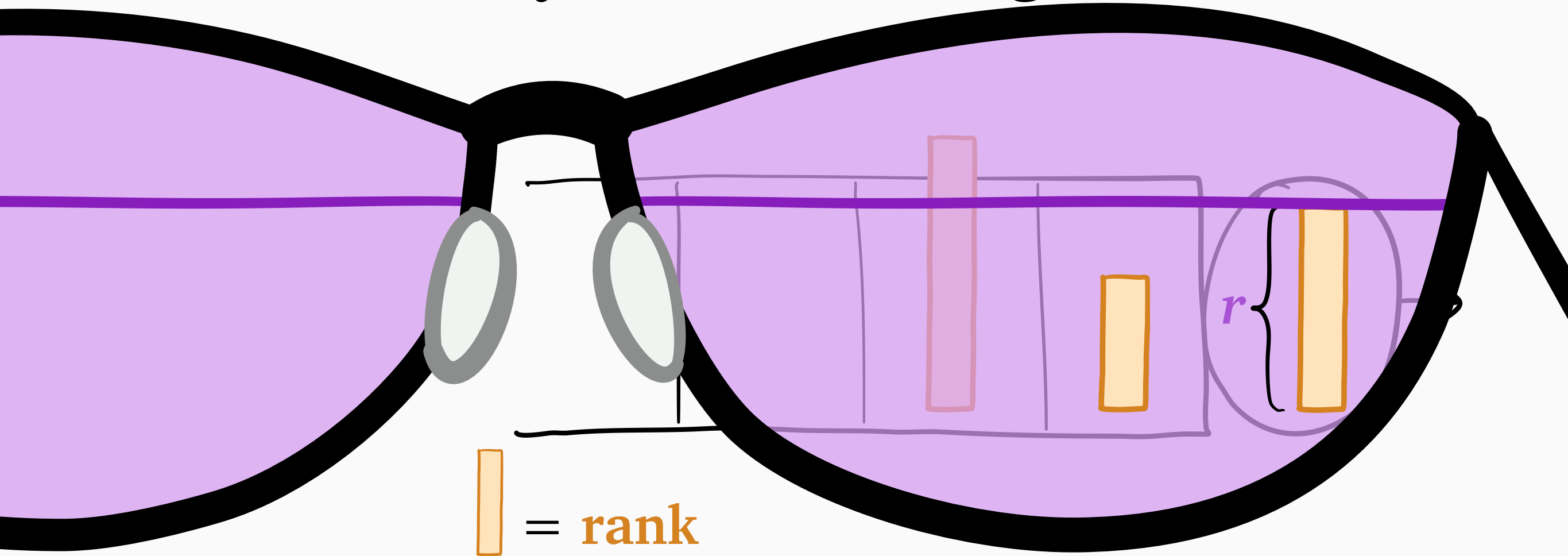
**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

# Key **SOAP** insight



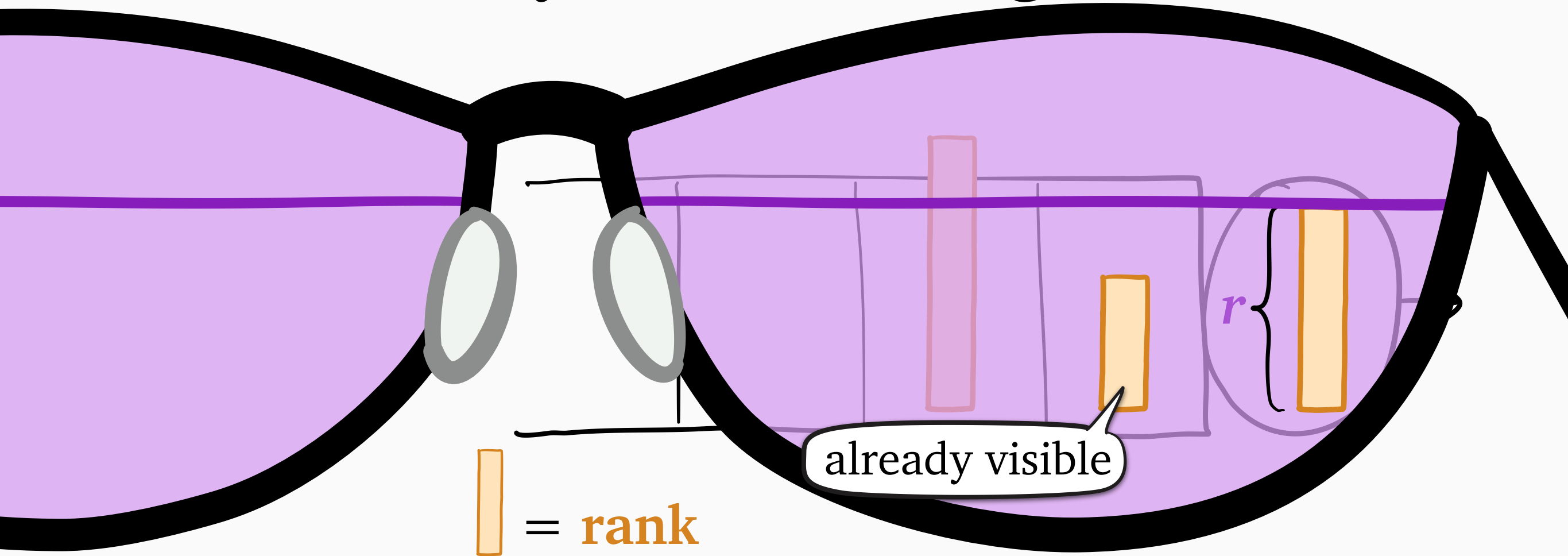
**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

# Key **SOAP** insight



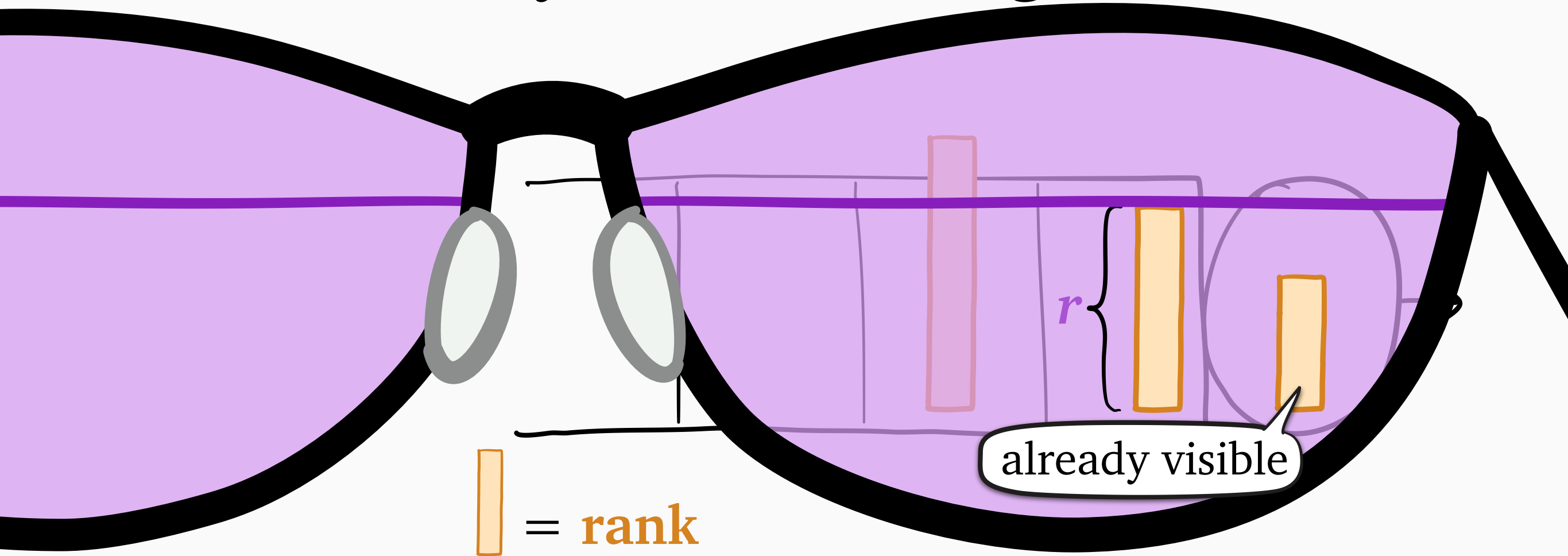
**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

# Key **SOAP** insight



**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

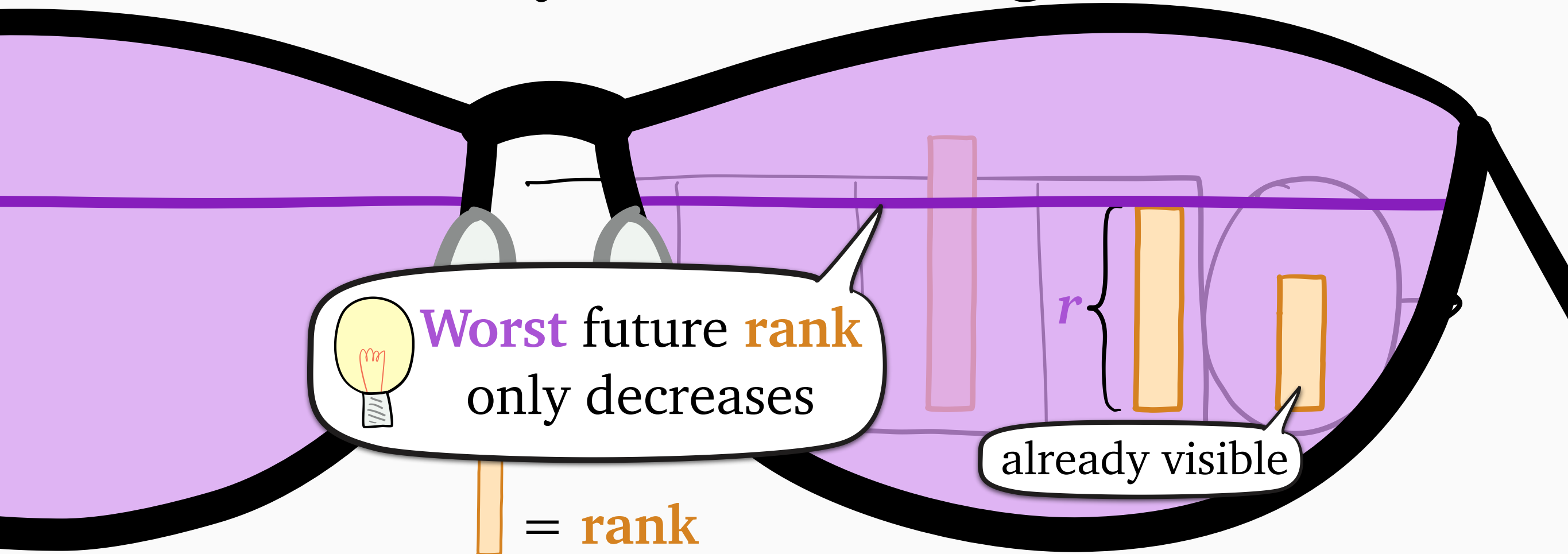
# Key **SOAP** insight



**Pessimism Principle:** compute  $W(r)$   
with  $r$  = my *worst future rank*

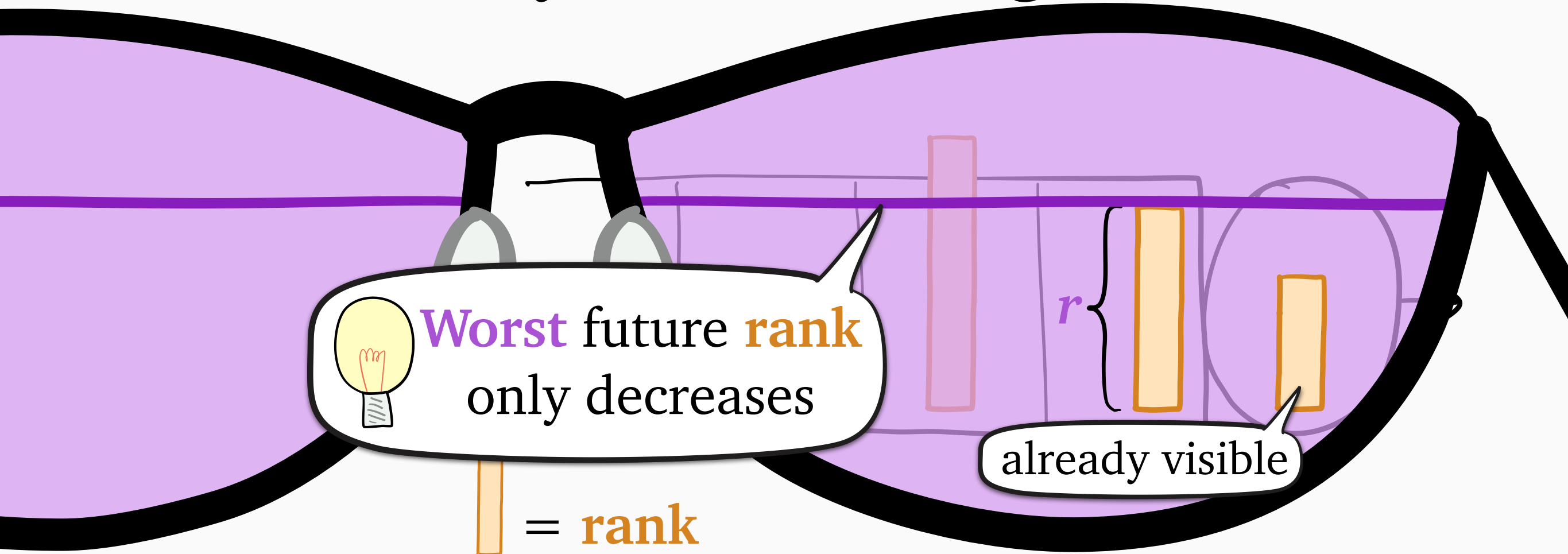


# Key SOAP insight



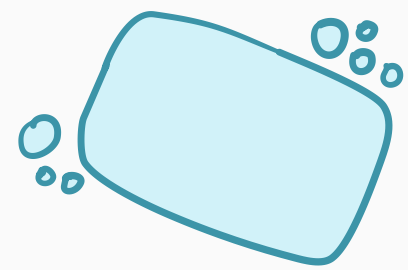
**Pessimism Principle:** compute  $W(r)$   
with  $r =$  my *worst future rank*

# Key SOAP insight



exact!

**Pessimism Principle:** compute  $W(r)$   
with  $r =$  my *worst future rank*



# Impact of SOAP



# Impact of SOAP

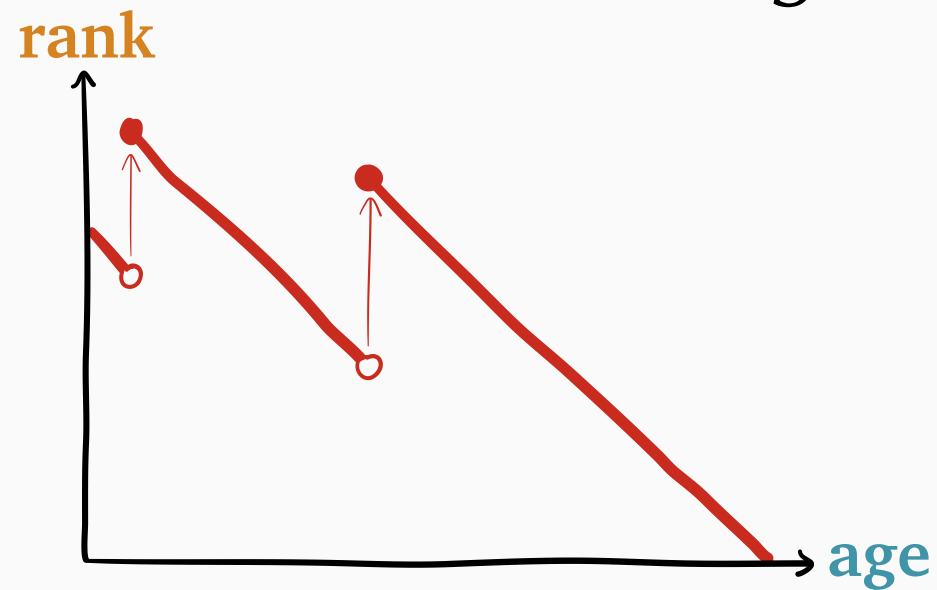


unknown sizes

# First analyses via SOAP

**SERPT**

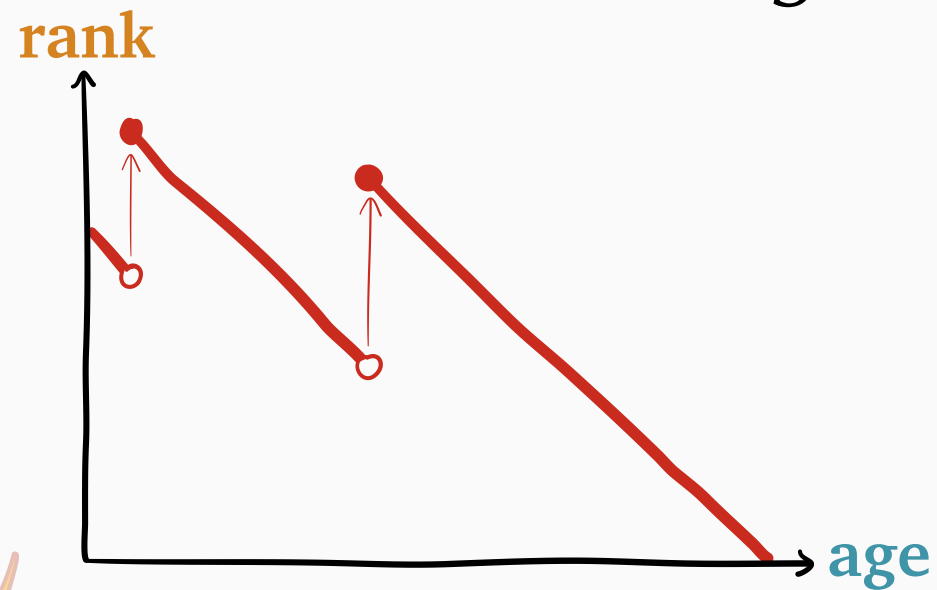
$$S = \begin{cases} 1 & \text{w.p. } \frac{1}{3} \\ 6 & \text{w.p. } \frac{1}{3} \\ 14 & \text{w.p. } \frac{1}{3} \end{cases}$$



# First analyses via SOAP

**SERPT**

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

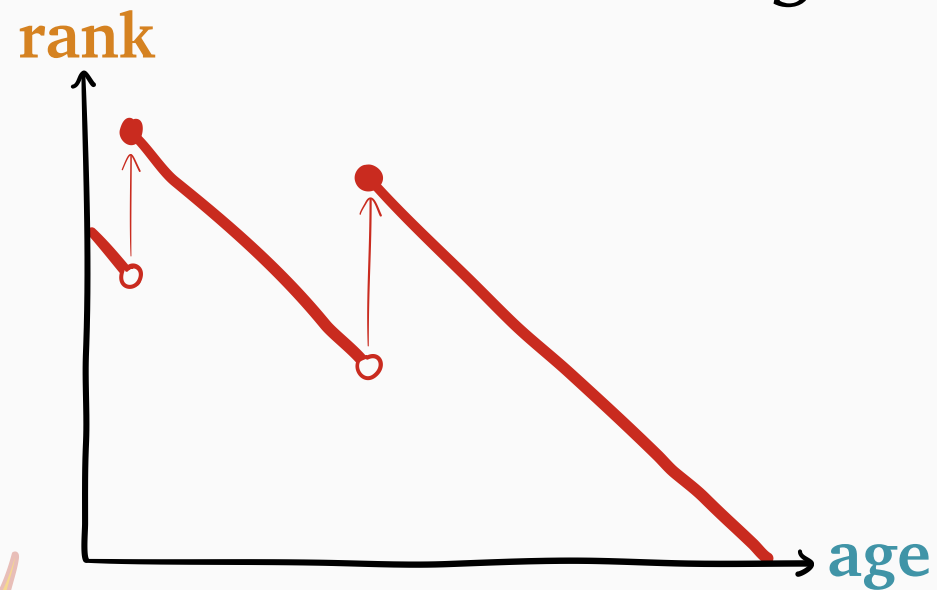


*first analysis of* **SERPT**

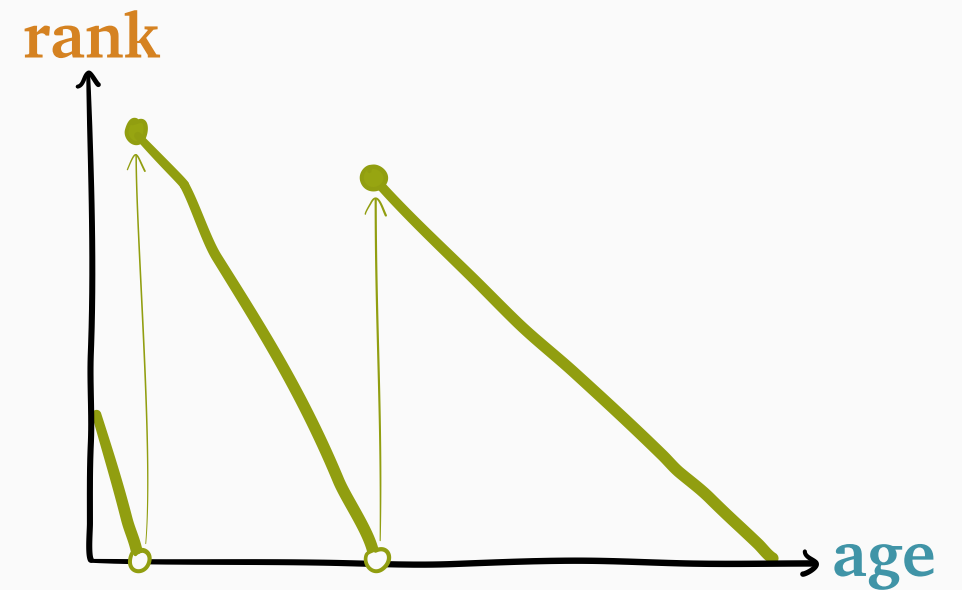
# First analyses via SOAP

**SERPT**

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



**Gittins**  
minimizes  $E[T]$

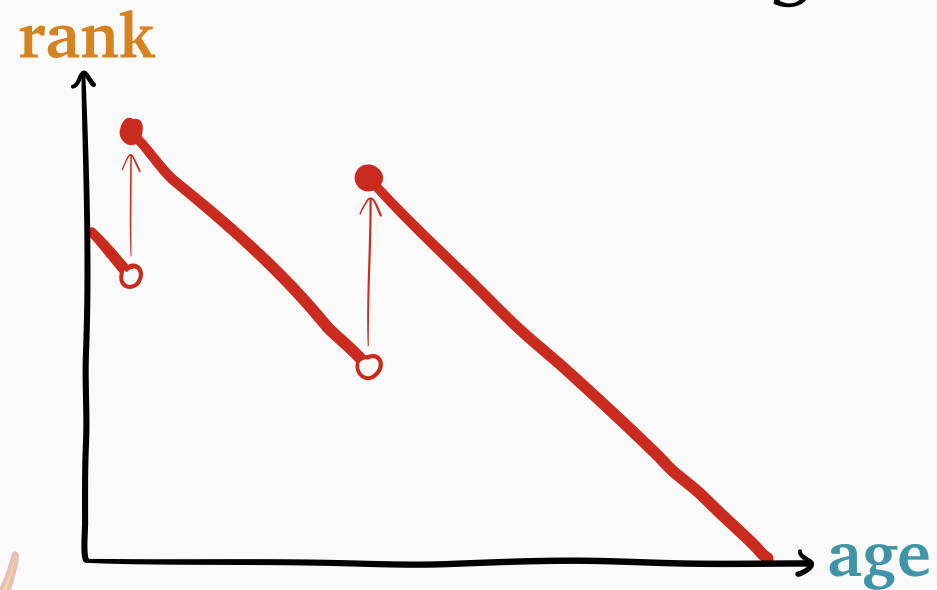


**NEW!** first analysis of **SERPT**

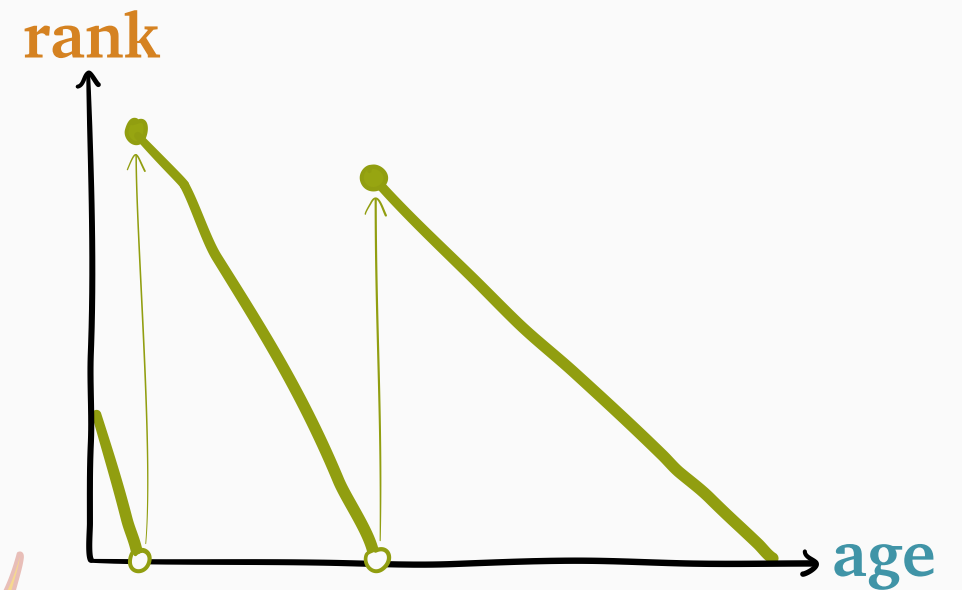
# First analyses via SOAP

**SERPT**

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



**Gittins**  
minimizes  $E[T]$



**NEW!** first analysis of **SERPT**

**NEW!** first analysis of **Gittins**

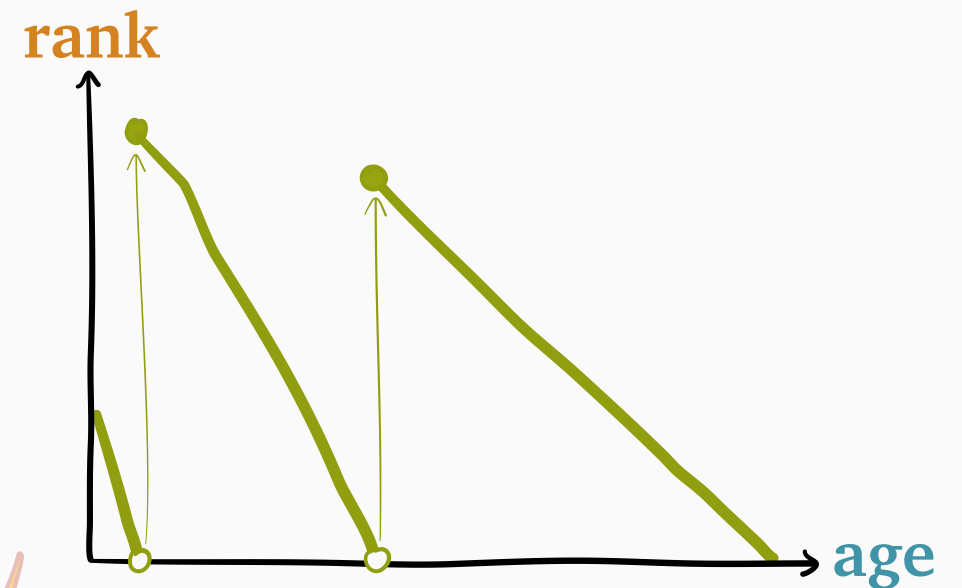
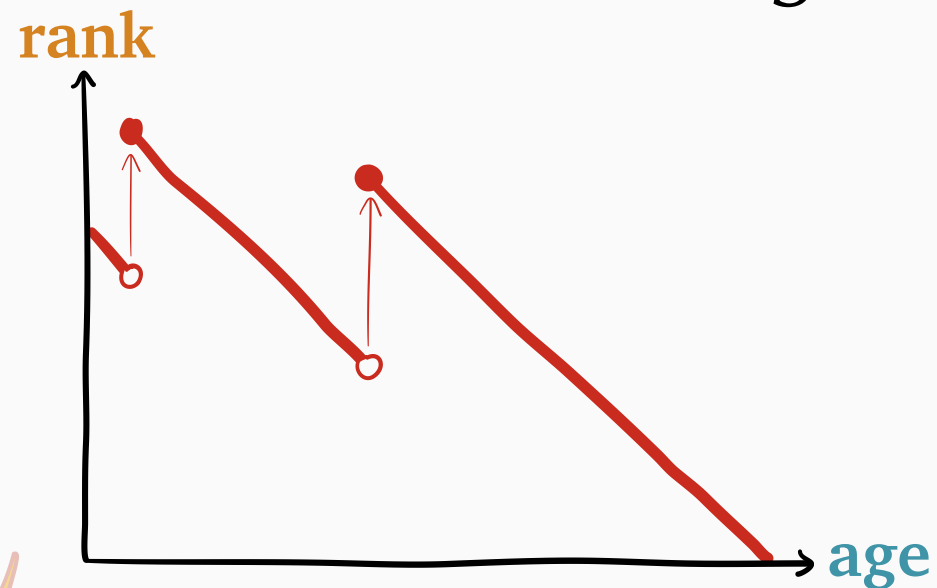


# First analyses via SOAP

**SERPT**

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

**Gittins**  
minimizes  $E[T]$



**NEW!**

first analysis of **SERPT**

**NEW!**

first analysis of **Gittins**

**NEW!**

Sample result: for heavy-tailed  $S$ , **SERPT** and **Gittins** have optimal tail decay in large- $t$  limit:

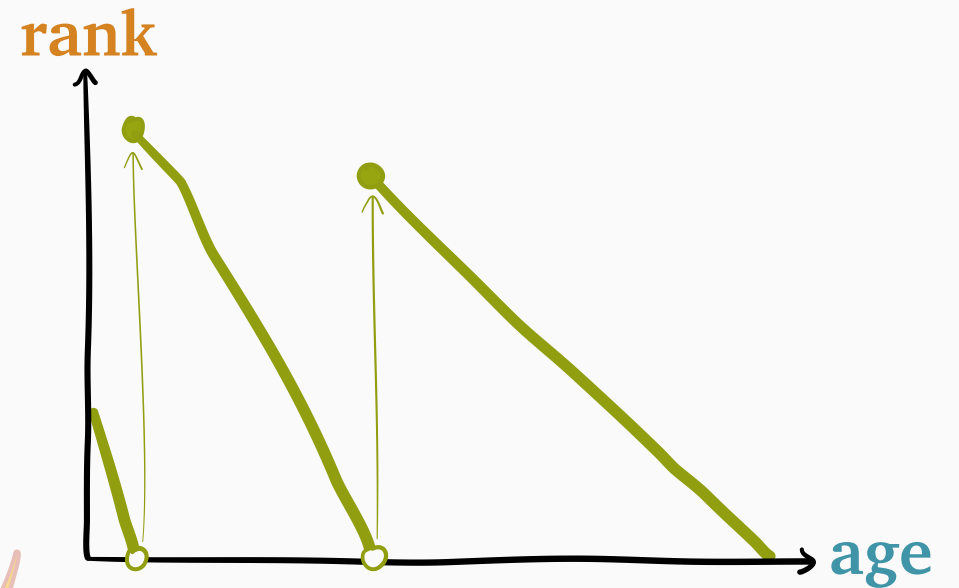
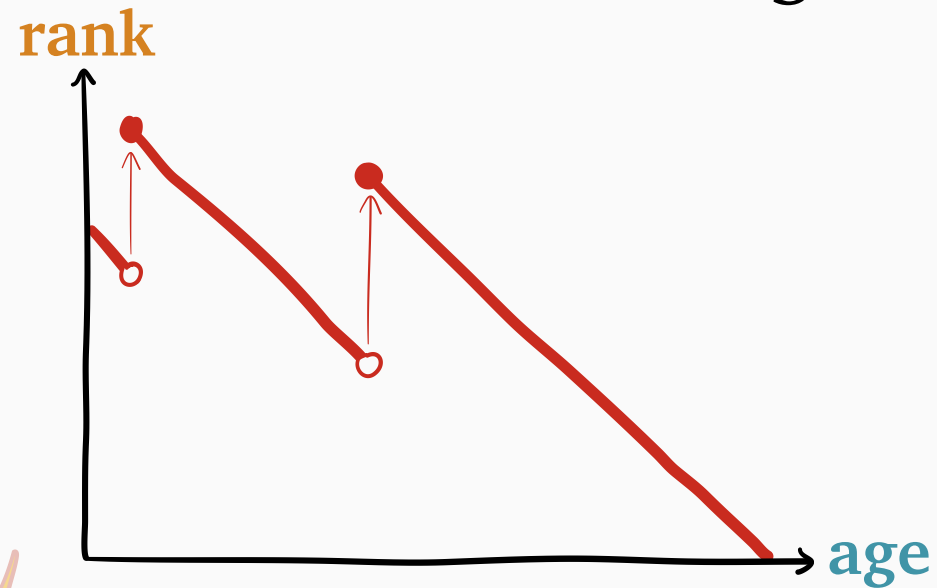
$$\mathbf{P}[T > t] \sim \mathbf{P}[S > (1 - \rho)t]$$

# First analyses via SOAP

**SERPT**

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

**Gittins**  
minimizes  $E[T]$



**NEW!** first analysis of **SERPT**

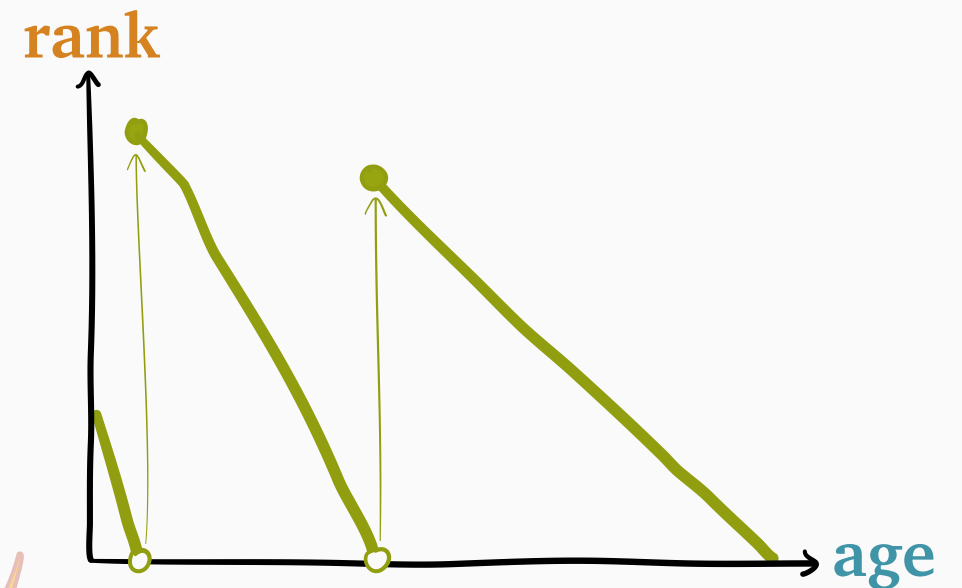
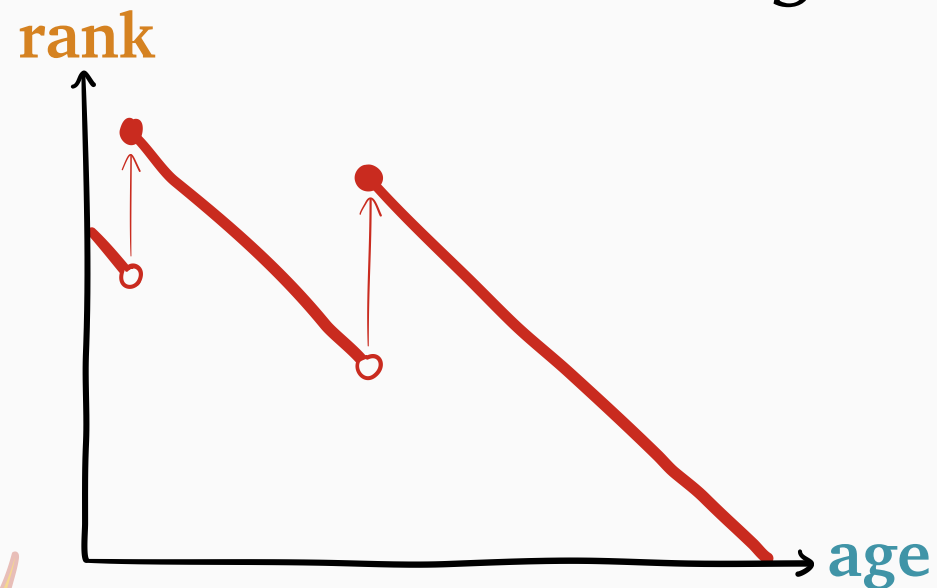
**NEW!** first analysis of **Gittins**

# First analyses via SOAP

**SERPT**

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

**Gittins**  
minimizes  $E[T]$



**NEW!** first analysis of **SERPT**

**NEW!** first analysis of **Gittins**

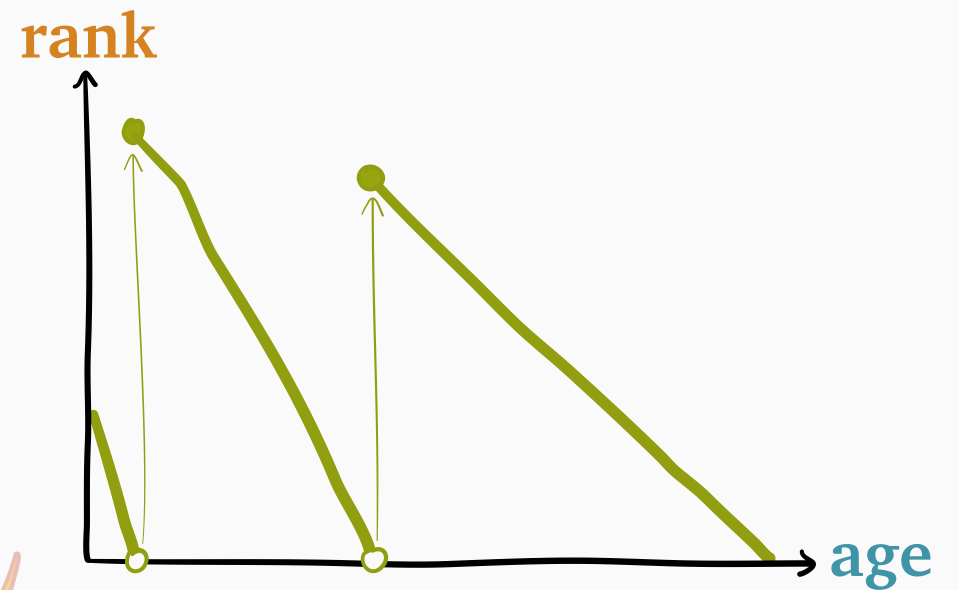
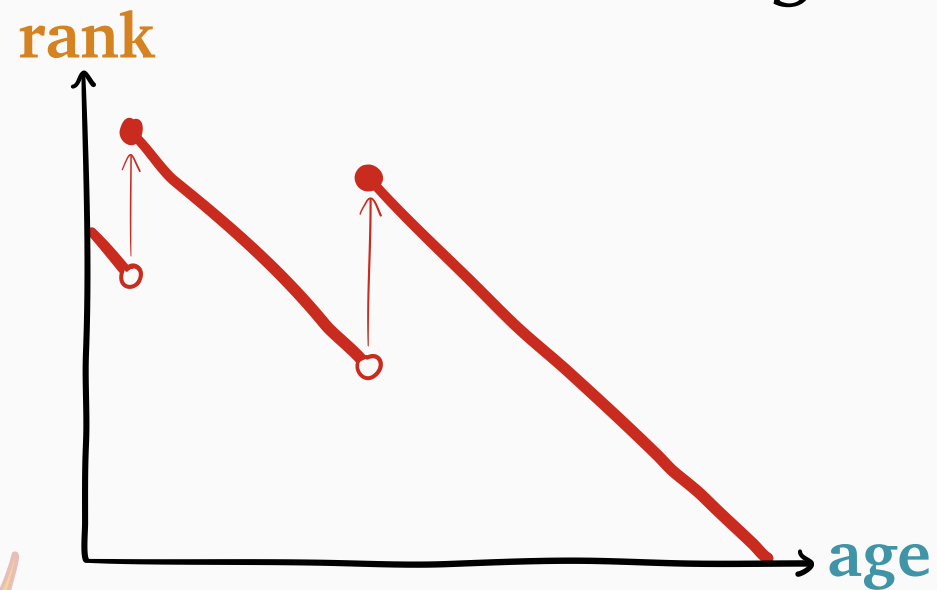
Observation: **SERPT**'s  $E[T]$  "usually" within 10% of **Gittins**'s

# First analyses via SOAP

**SERPT**

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

**Gittins**  
minimizes  $E[T]$



**NEW!** first analysis of **SERPT**

**NEW!** first analysis of **Gittins**

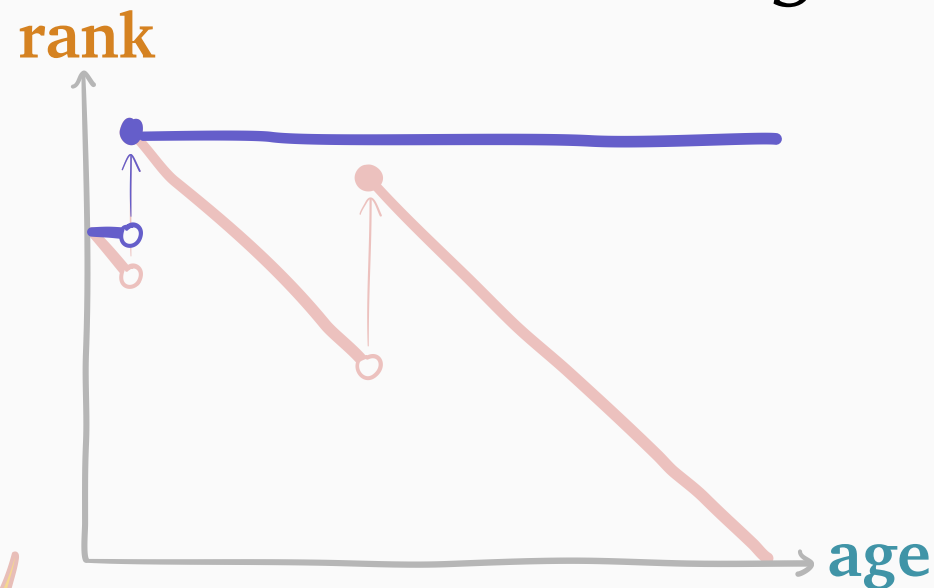
Observation: **SERPT**'s  $E[T]$  "usually" within 10% of **Gittins**'s

 no proof yet

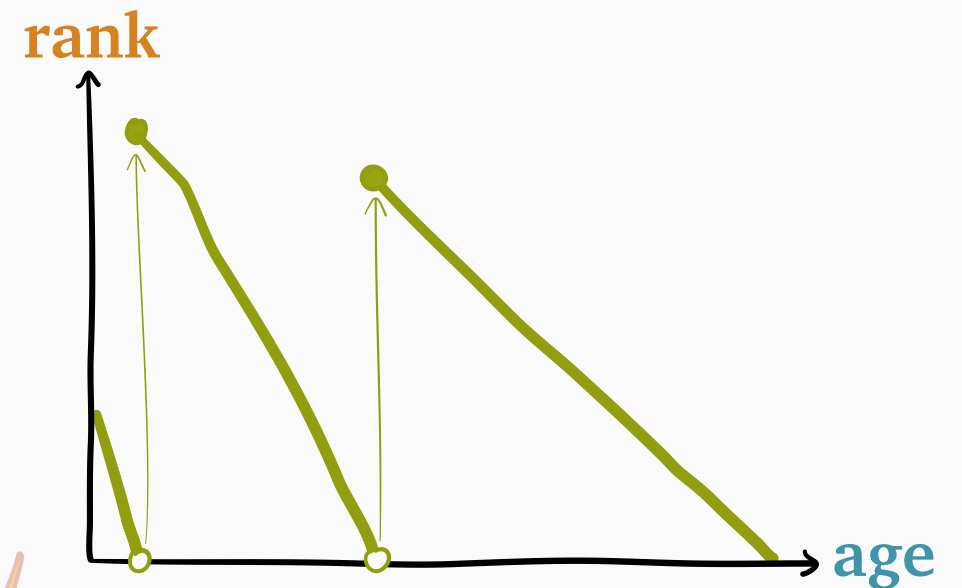
# First analyses via SOAP

M-SERP

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



Gittins  
minimizes  $E[T]$



**NEW!** first analysis of **SERP**

**NEW!** first analysis of **Gittins**

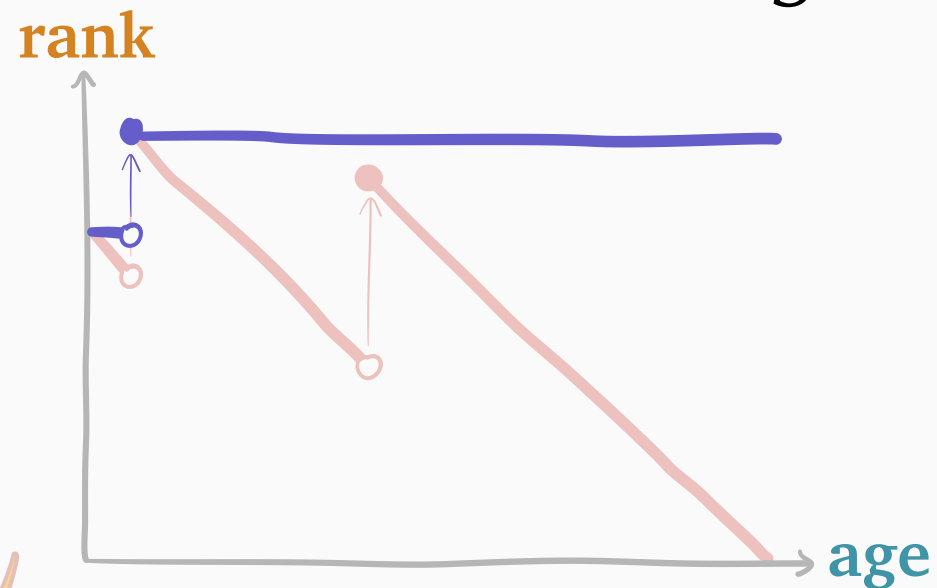
Observation: **SERP**'s  $E[T]$  "usually" within 10% of **Gittins**'s

 no proof yet

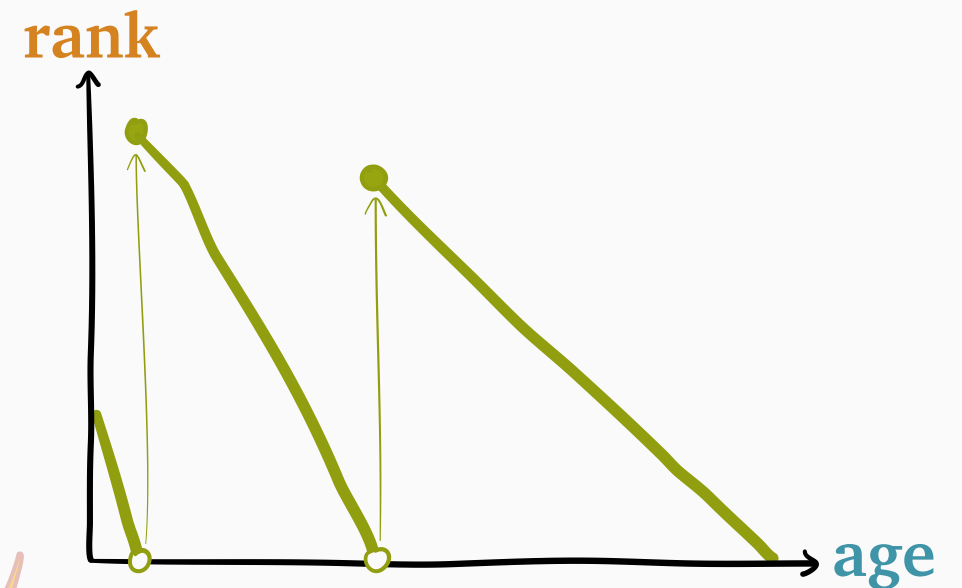
# First analyses via SOAP

**M-SERPT**

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



**Gittins**  
minimizes  $E[T]$



**NEW!** first analysis of **SERPT**

**NEW!** first analysis of **Gittins**

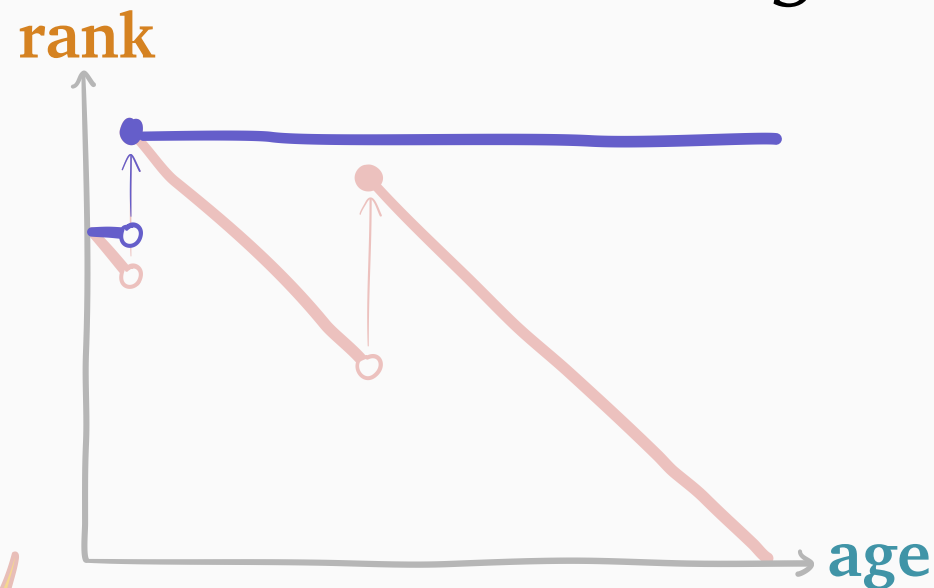
**Observation:** **SERPT**'s  $E[T]$  “usually” within 10% of **Gittins**'s

**Theorem:** **M-SERPT** is a 5-approximation for  $E[T]$

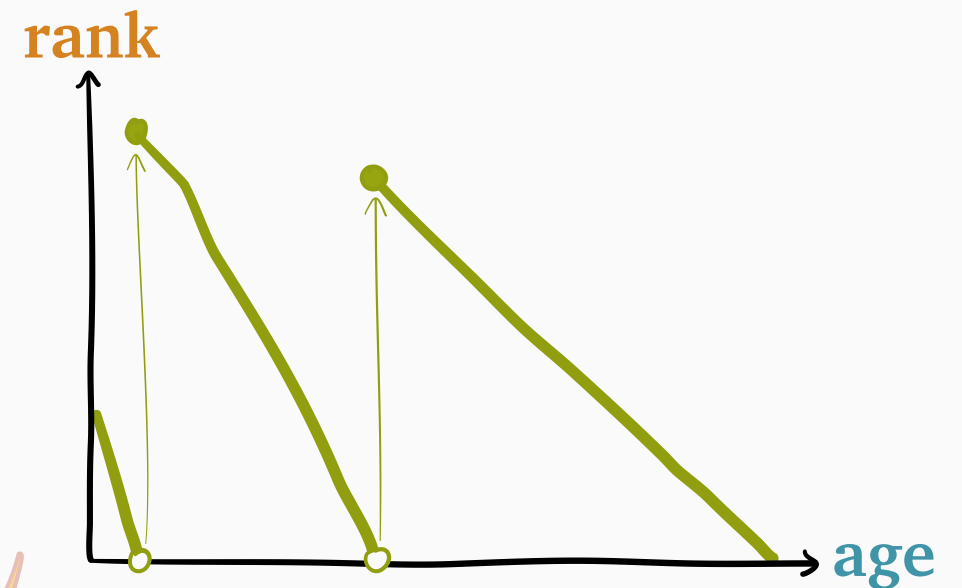
# First analyses via SOAP

**M-SERPT**

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



**Gittins**  
minimizes  $E[T]$



first analysis of **SERPT**



first analysis of **Gittins**

**Lesson:** carefully plan how **rank** gets worse

Observed values of  $E[T]$  "usually" within 10% of **Gittins's**

**Theorem:** **M-SERPT** is a 5-approximation for  $E[T]$



# Impact of SOAP



unknown sizes





# Impact of SOAP



unknown sizes



*First analysis of*

- mean
- tail decay



# Impact of SOAP



unknown sizes



preemption limitations



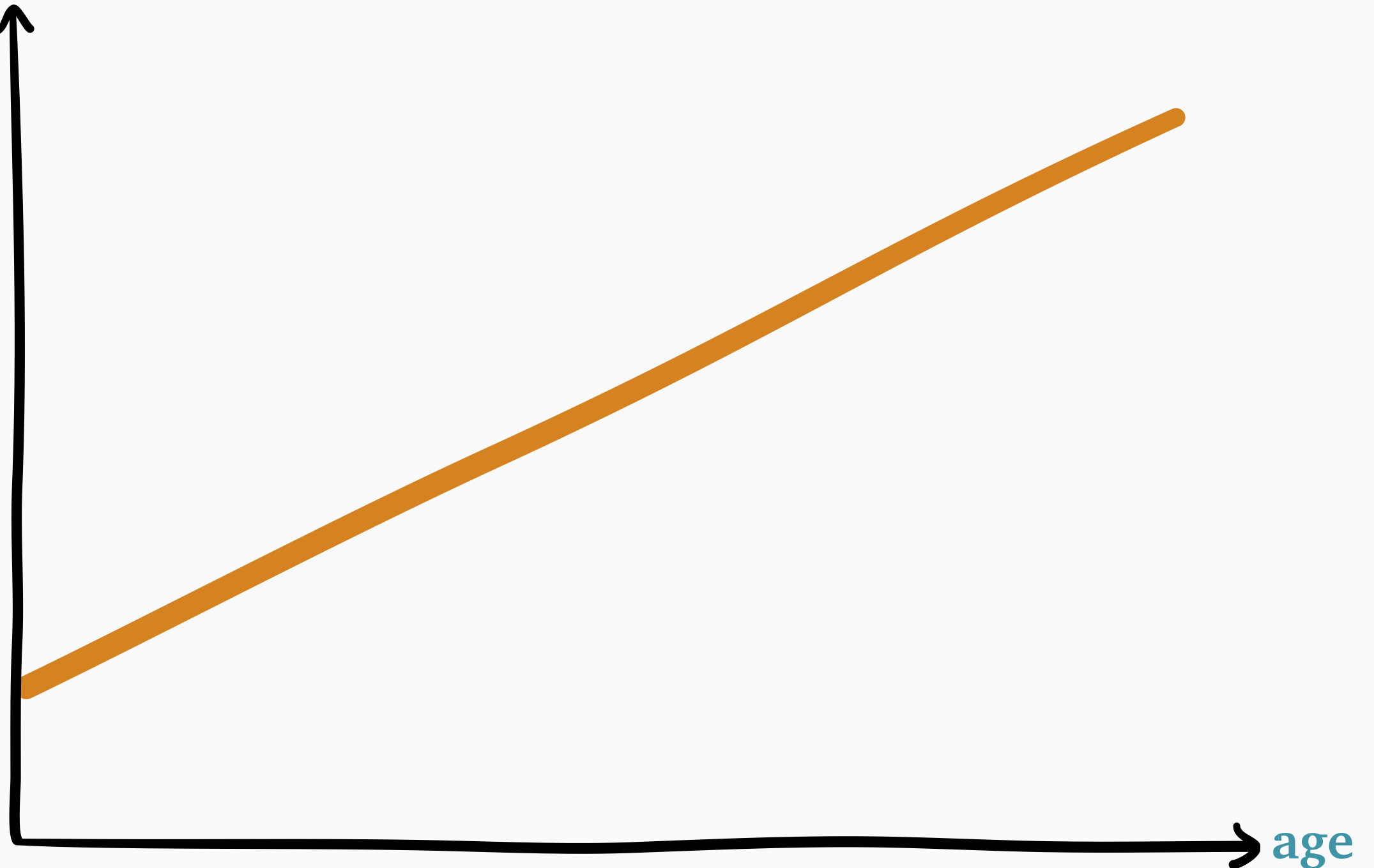
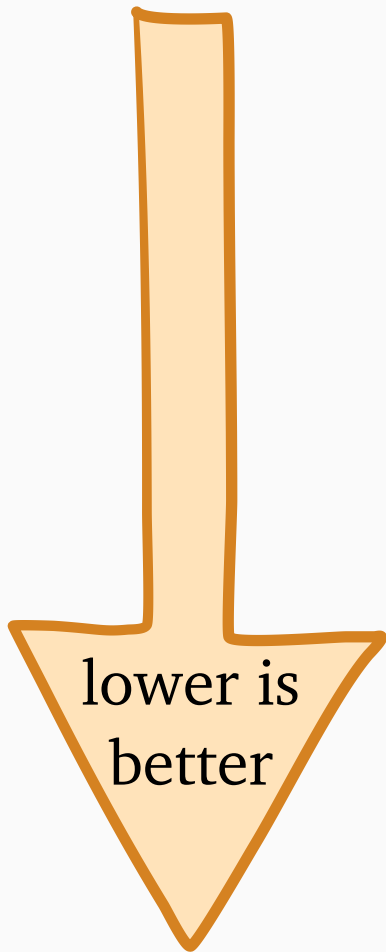
*First analysis of*

- mean
- tail decay

# Preemption checkpoints

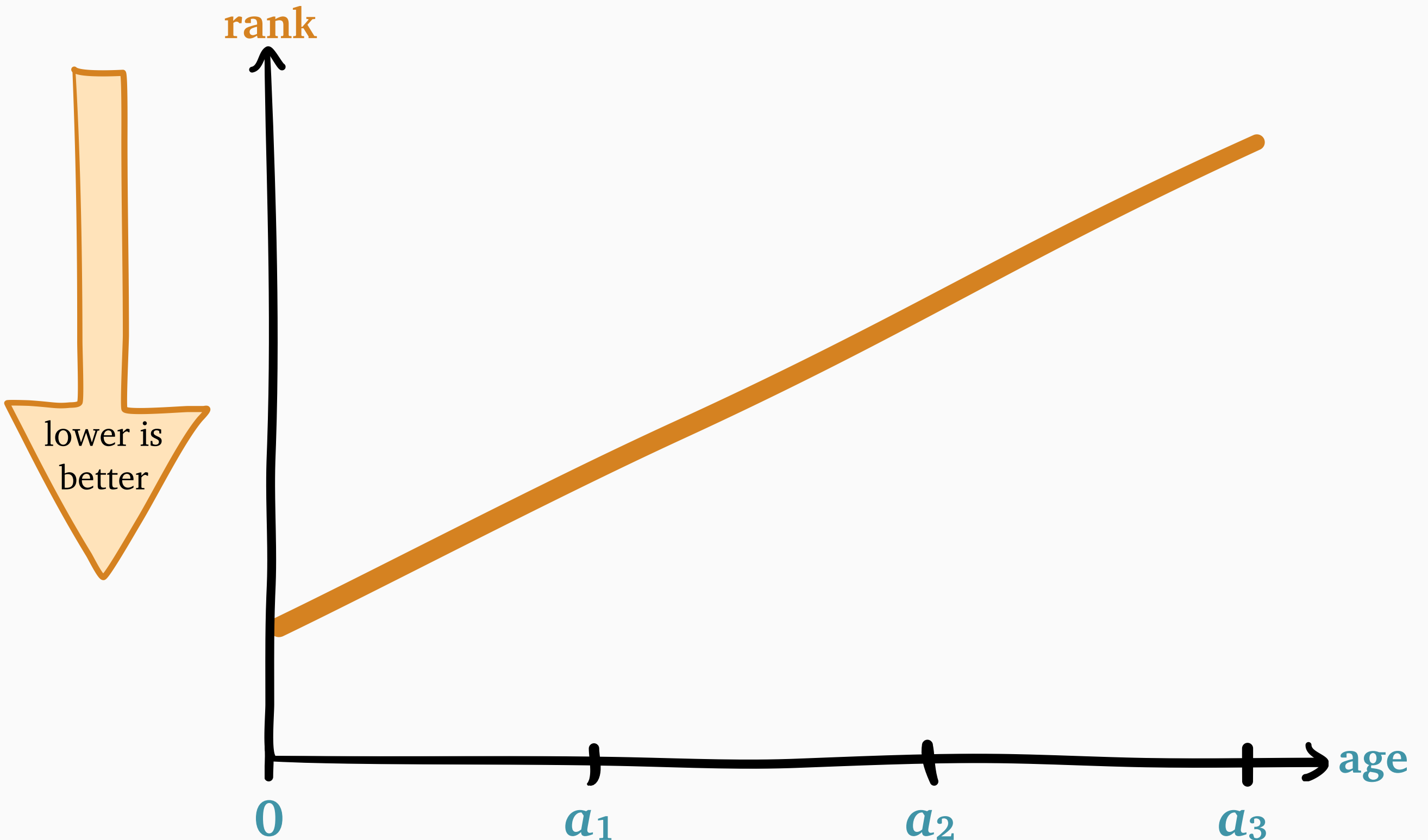
Can only preempt only at *checkpoint ages*  $a_i$

rank



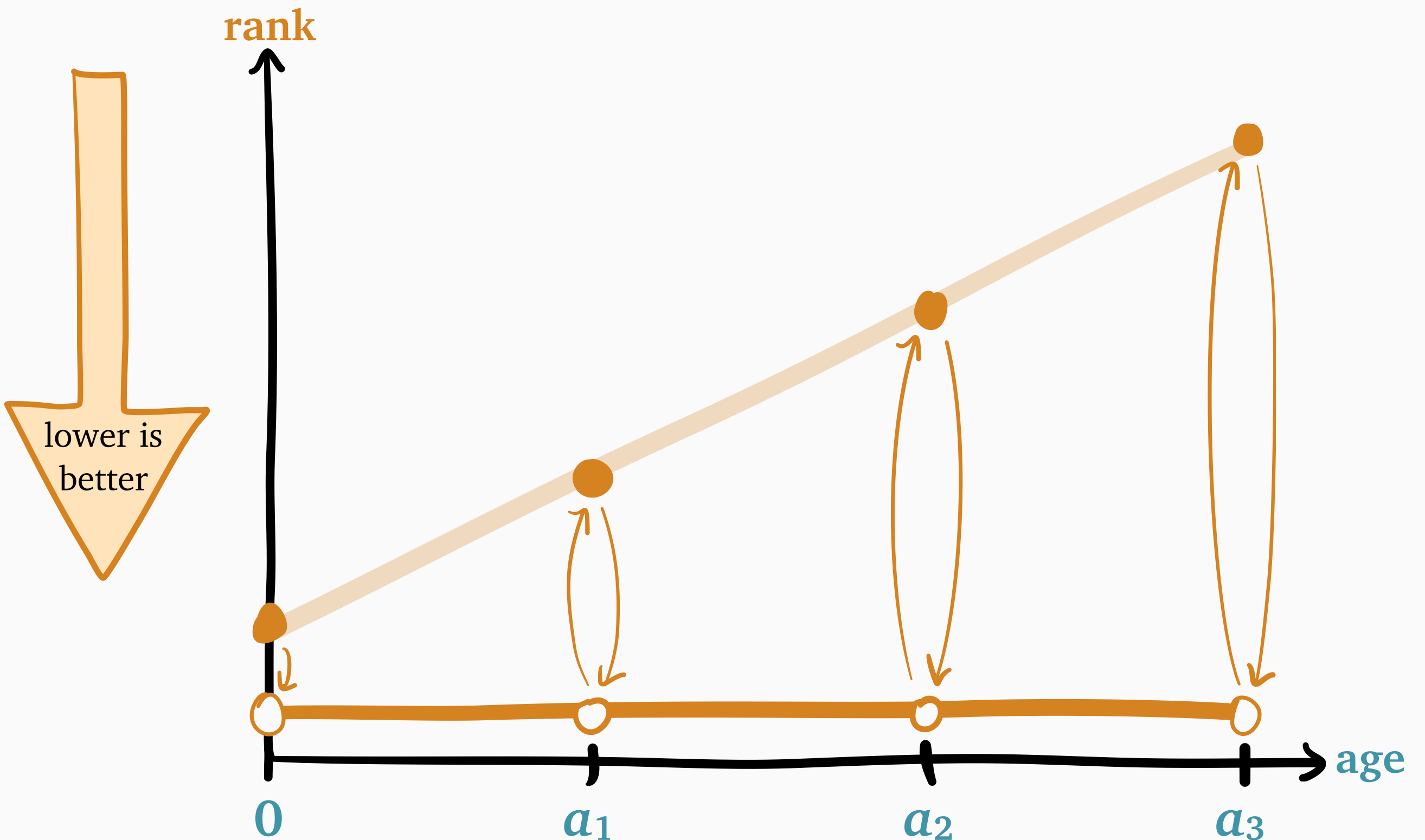
# Preemption checkpoints

Can only preempt only at *checkpoint ages*  $a_i$



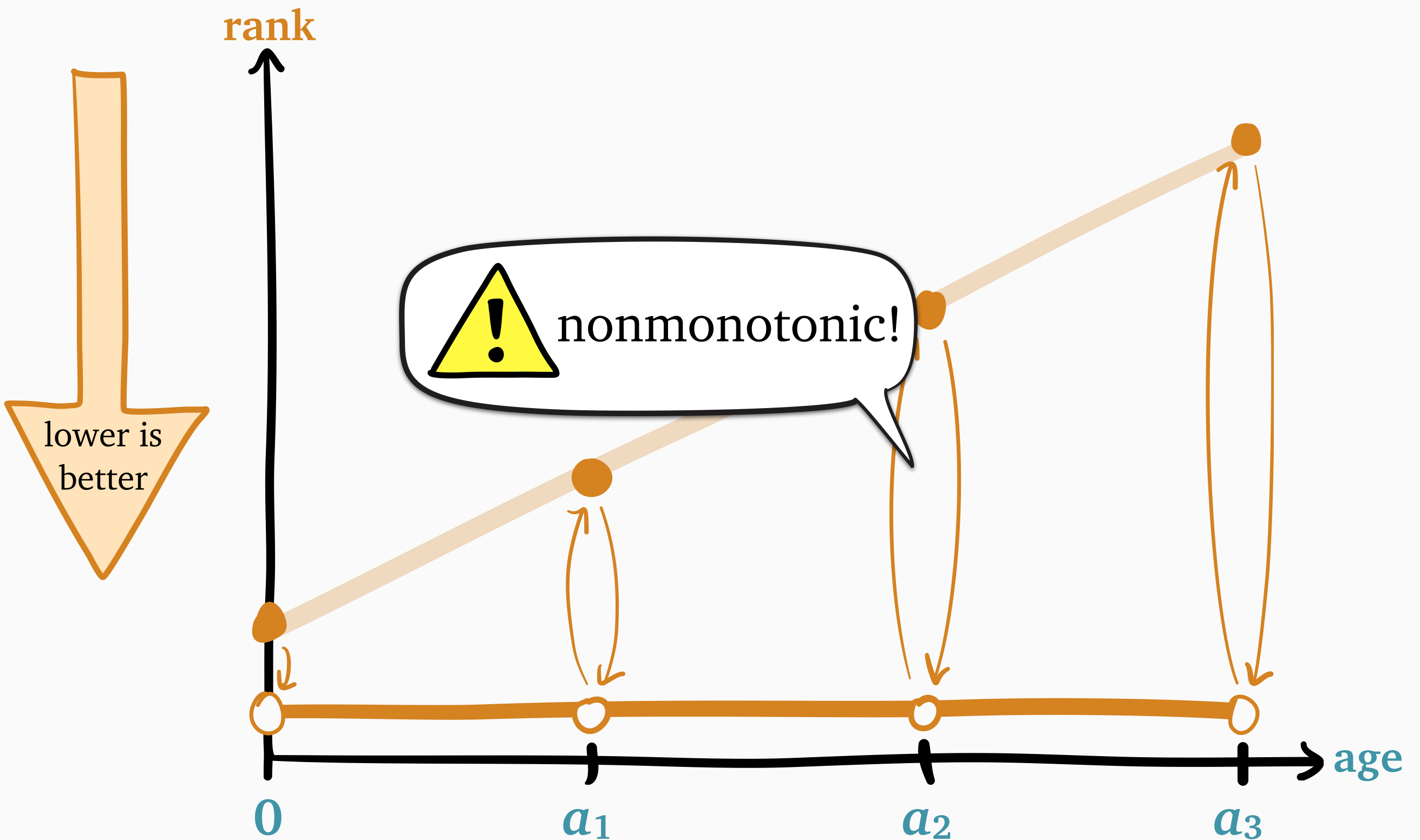
# Preemption checkpoints

Can only preempt only at *checkpoint ages*  $a_i$



# Preemption checkpoints

Can only preempt only at *checkpoint ages*  $a_i$





# Impact of SOAP



unknown sizes

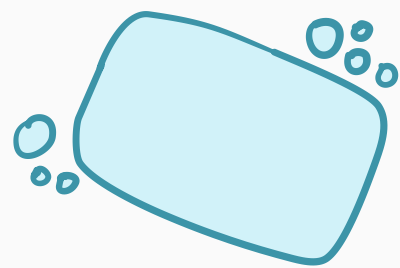


preemption limitations



*First analysis of*

- mean
- tail decay



# Impact of SOAP



unknown sizes



preemption limitations



*First analysis of*

- mean
- tail decay



*First analysis of*

- mean
- tail decay



# SOAP References

- [1] Scully, Harchol-Balter, and Scheller-Wolf (2018). “SOAP: One Clean Analysis of All Age-Based Scheduling Policies.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018). **Finalist: 2019 INFORMS APS Best Student Paper Prize.**
- [2] Scully and Harchol-Balter (2018). “SOAP Bubbles: Robust Scheduling Under Adversarial Noise.” *56th Annual Allerton Conference on Communication, Control, and Computing.*
- [3] Scully, Harchol-Balter, and Scheller-Wolf (2020). “Simple Near-Optimal Scheduling for the M/G/1.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018). **Winner: SIGMETRICS 2020 Best Video Award.**
- [4] Scully, Keveld, Boxma, Dorsman, and Wierman (2020). “Characterizing policies with optimal response time tails under heavy-tailed job sizes.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018).
- [5] Scully and Kreveld (2021). “When Does the Gittins Policy Have Asymptotically Optimal Response Time Tail?” *SIGMETRICS Perform. Eval. Rev.* (MAMA 2021), *arXiv*.
- [6] Scully and Harchol-Balter (2021). “How to Schedule Near-Optimally under Real-World Constraints.” *arXiv*.

# SOAP References

[1] Scully, Harchol-Balter, and Scheller-Wolf (2018). “SOAP: One Clean Analysis of All Age-Based Scheduling Policies.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018). **Finalist: 2019 INFORMS APS Best Student Paper Prize.**

[2] Scully and Harchol-Balter (2018). “SOAP Bubbles: Robust Scheduling Under Adversarial Noise.” *56th Annual Allerton Conference on Communication, Control, and Computing.*

[3] Scully, Harchol-Balter, and Scheller-Wolf (2020). “Simple Near-Optimal Scheduling for the M/G/1.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018). **Winner: SIGMETRICS 2020 Best Video Award.**

[4] Scully, Keveld, Boxma, Dorsman, and Wierman (2020). “Characterizing policies with optimal response time tails under heavy-tailed job sizes.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018).

[5] Scully and Kreveld (2021). “When Does the Gittins Policy Have Asymptotically Optimal Response Time Tail?” *SIGMETRICS Perform. Eval. Rev.* (MAMA 2021), *arXiv*.

[6] Scully and Harchol-Balter (2021). “How to Schedule Near-Optimally under Real-World Constraints.” *arXiv*.



# SOAP References

[1] Scully, Harchol-Balter, and Scheller-Wolf (2018). “SOAP: One Clean Analysis of All Age-Based Scheduling Policies.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018). **Finalist: 2019 INFORMS APS Best Student Paper Prize.**

[2] Scully and Harchol-Balter (2018). “SOAP Bubbles: Robust Scheduling Under Adversarial Noise.” *56th Annual Allerton Conference on Communication, Control, and Computing.*

[3] Scully, Harchol-Balter, and Scheller-Wolf (2020). “Simple Near-Optimal Scheduling for the M/G/1.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018). **Winner: SIGMETRICS 2020 Best Video Award.**

[4] Scully, Keveld, Boxma, Dorsman, and Wierman (2020). “Characterizing policies with optimal response time tails under heavy-tailed job sizes.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2018).

[5] Scully and Kreveld (2021). “When Does the Gittins Policy Have Asymptotically Optimal Response Time Tail?” *SIGMETRICS Perform. Eval. Rev.* (MAMA 2021), *arXiv*.

[6] Scully and Harchol-Balter (2021). “How to Schedule Near-Optimally under Real-World Constraints.” *arXiv*.

case  
studies

# Outline: two new tools



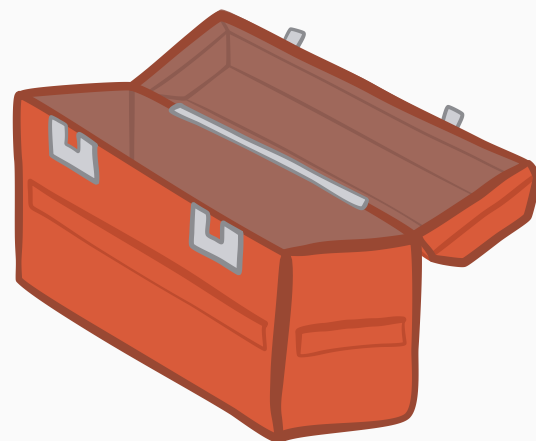
new unifying theory of  
single-server scheduling

- unknown sizes
- preemption limitations



new queueing identity to  
complement Little's Law

- multiserver systems
- noisy size estimates

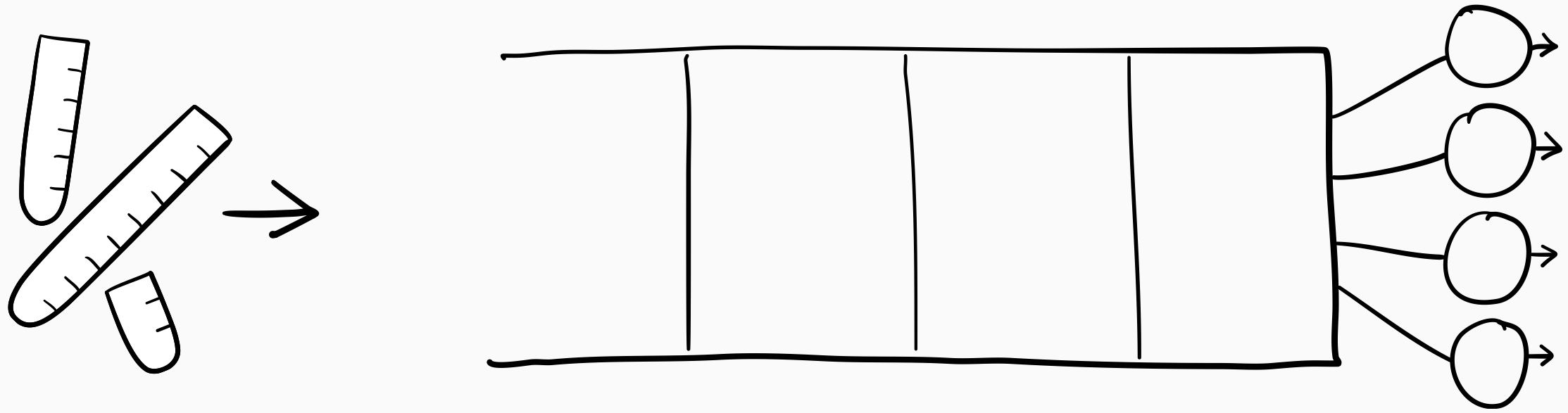




*Part 2*

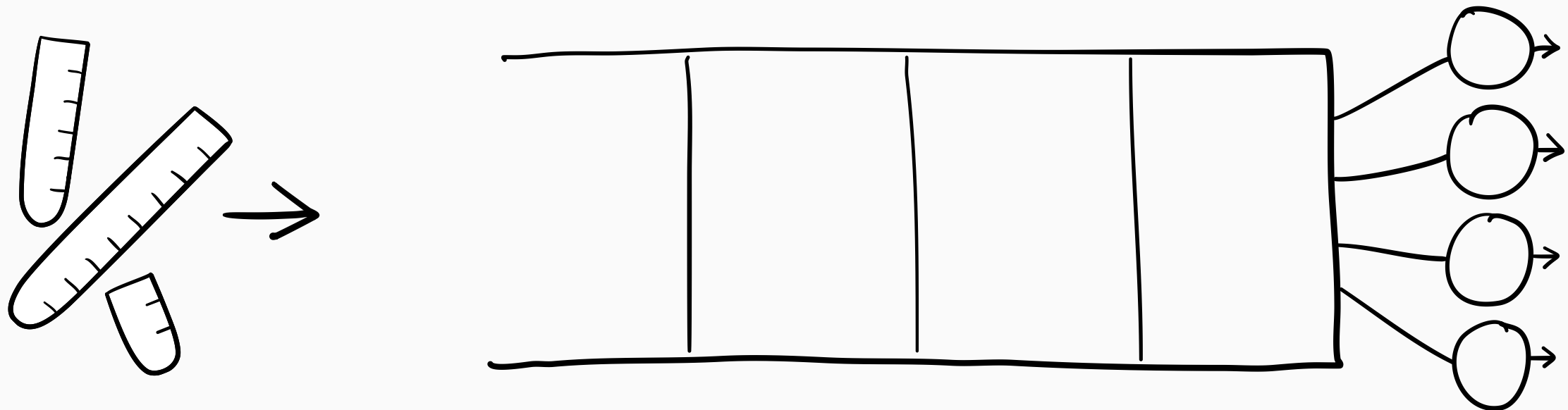
**WINE**

# Multiserver queueing system



# Multiserver queueing system

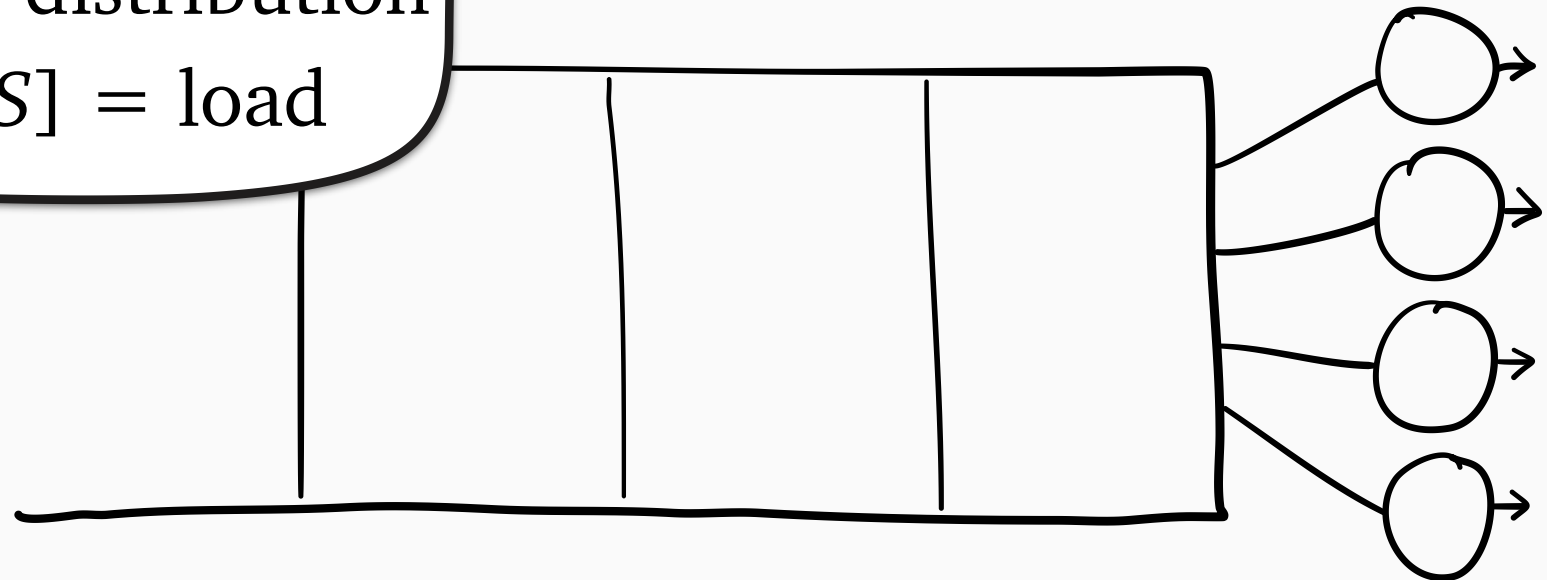
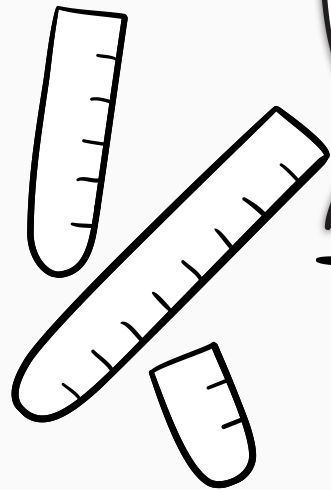
M/G/k



# Multiserver queueing system

M/G/k

$\lambda$  = arrival rate  
 $S$  = size distribution  
 $\rho = \lambda \mathbf{E}[S] = \text{load}$

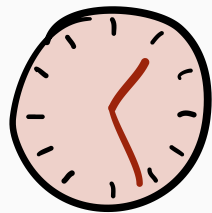
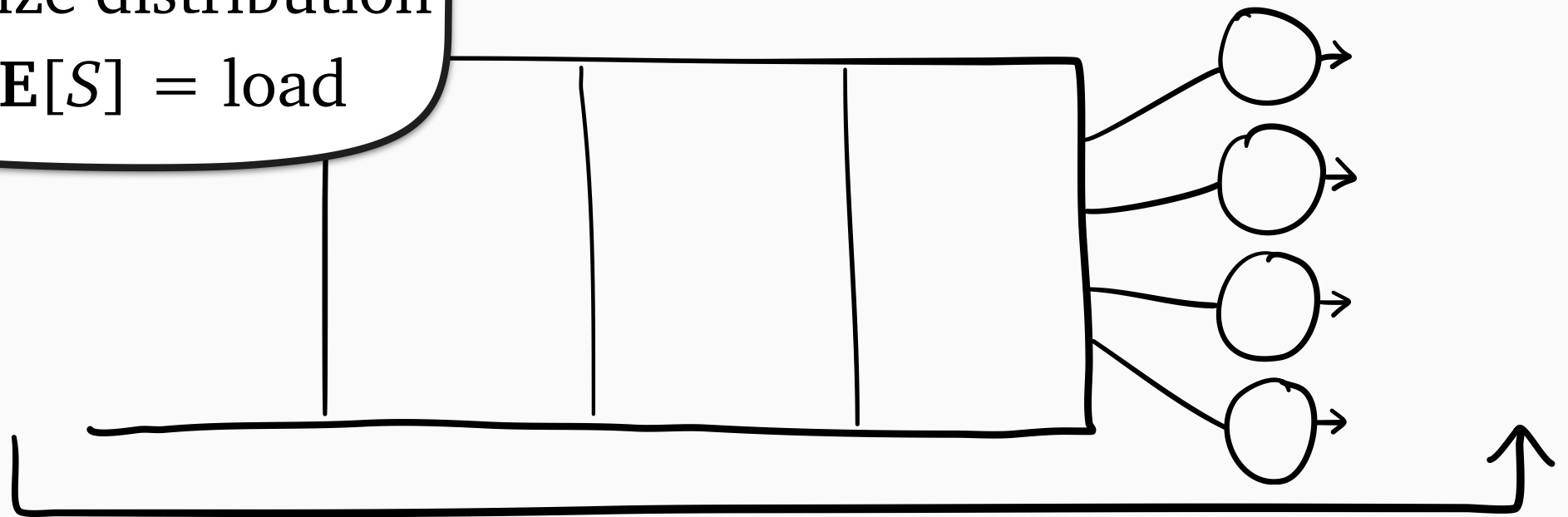
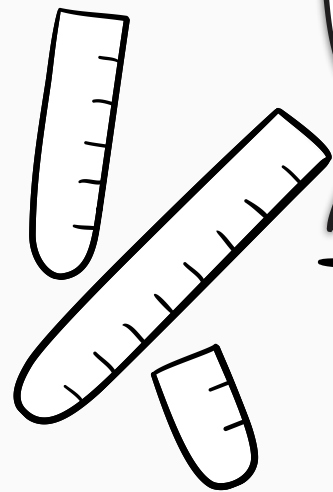




# Multiserver queueing system

M/G/k

$\lambda$  = arrival rate  
 $S$  = size distribution  
 $\rho = \lambda \mathbf{E}[S] = \text{load}$



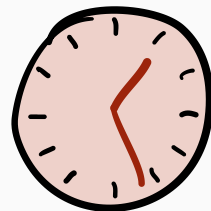
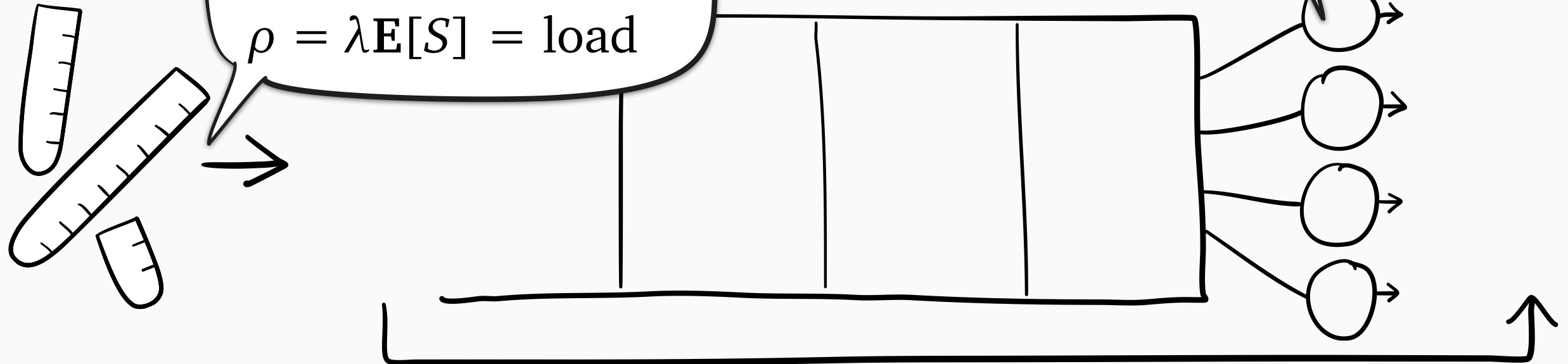
response time distribution  $T$

# Multiserver queueing system

M/G/*k*

$\lambda$  = arrival rate  
 $S$  = size distribution  
 $\rho = \lambda \mathbf{E}[S] = \text{load}$

*k* servers

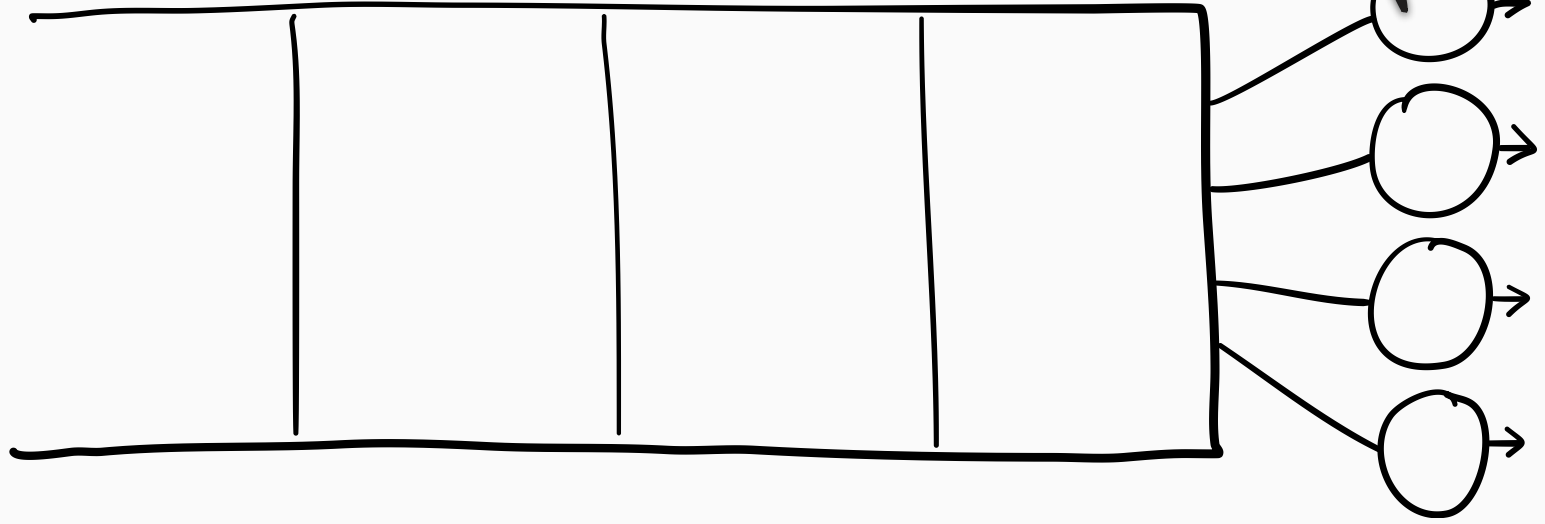
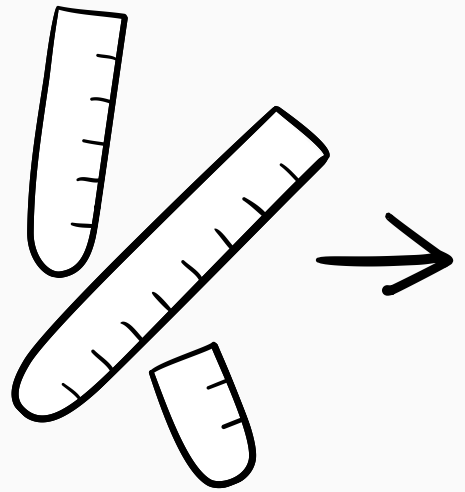


response time distribution  $T$

# Multiserver queueing system

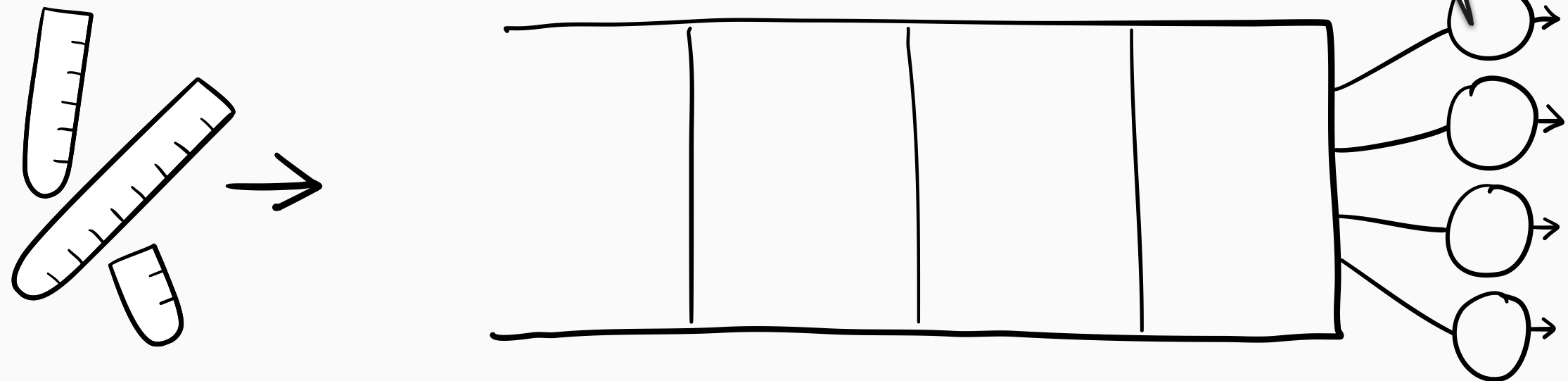
M/G/*k*

*k* servers



# Multiserver queueing system

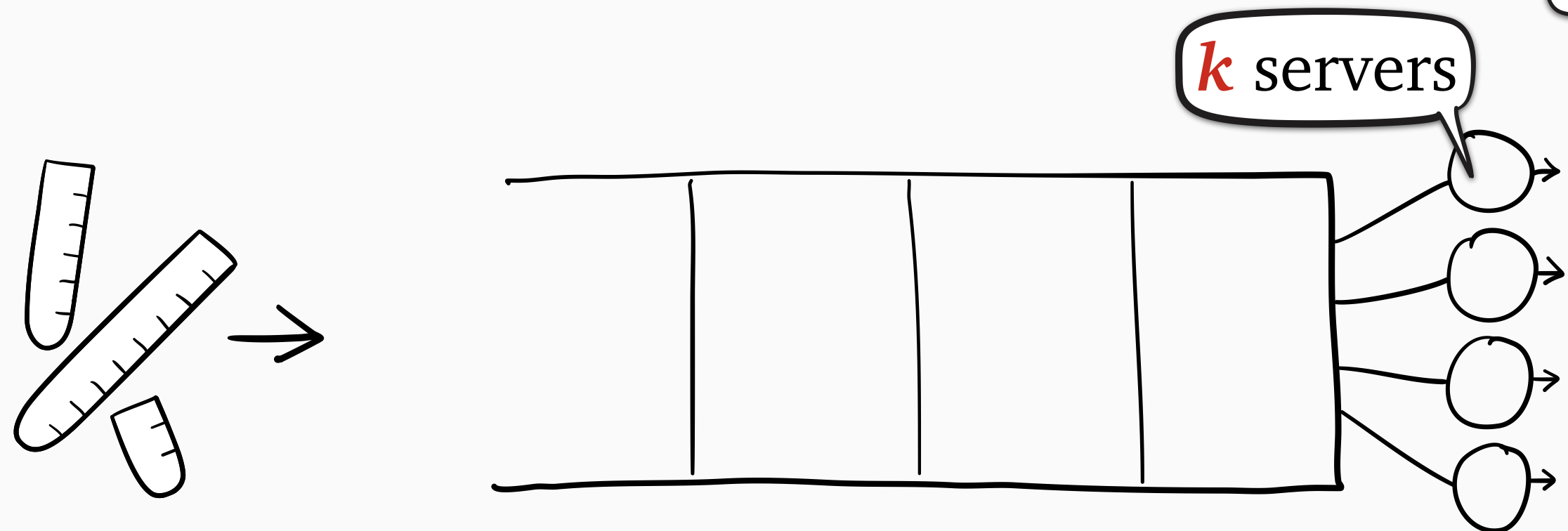
M/G/*k*



SRPT: **rank** = remaining size

# Multiserver queueing system

M/G/*k*

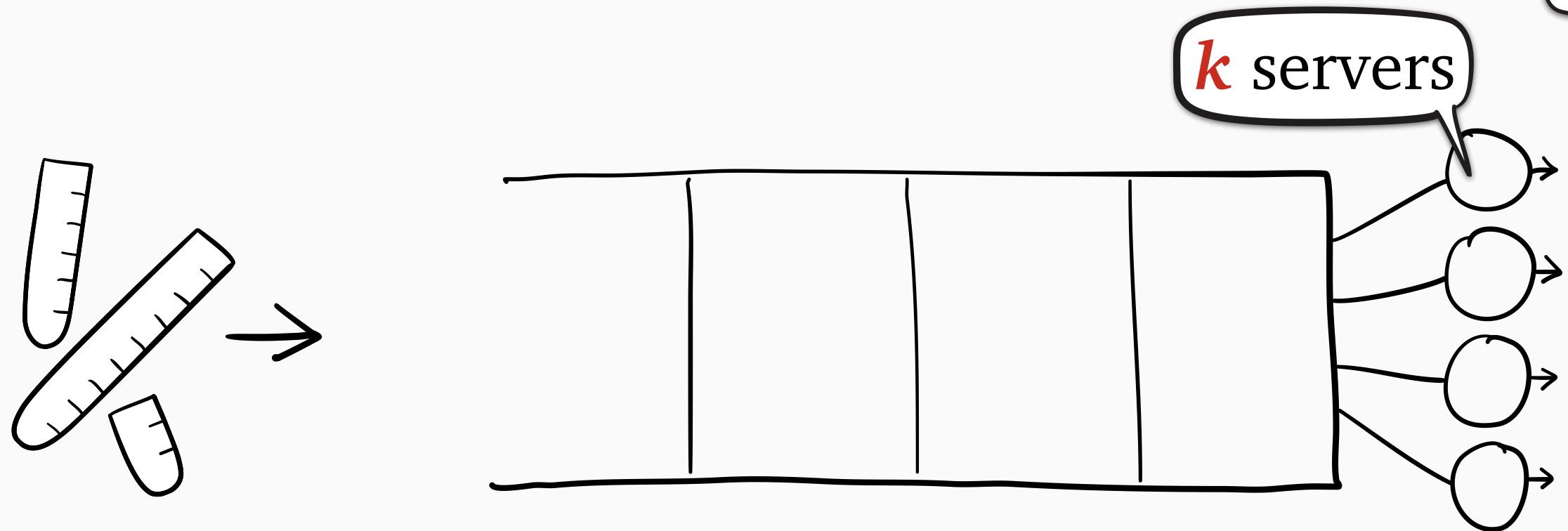


SRPT-1 (**single-server**): serves the job of least **rank**

SRPT: **rank** = remaining size

# Multiserver queueing system

M/G/ $k$



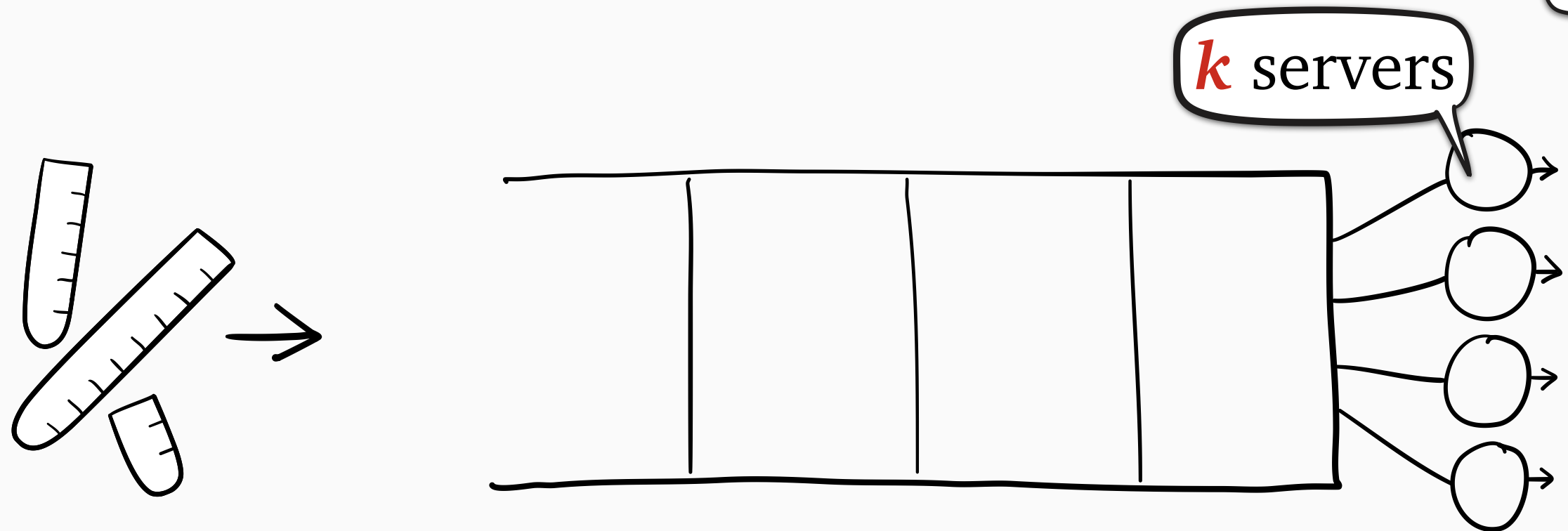
**SRPT-1 (single-server)**: serves the job of least **rank**

SRPT: **rank** = remaining size

**SRPT- $k$  (multiserver)**: serves the  $k$  jobs of least **rank**

# Multiserver queueing system

M/G/*k*



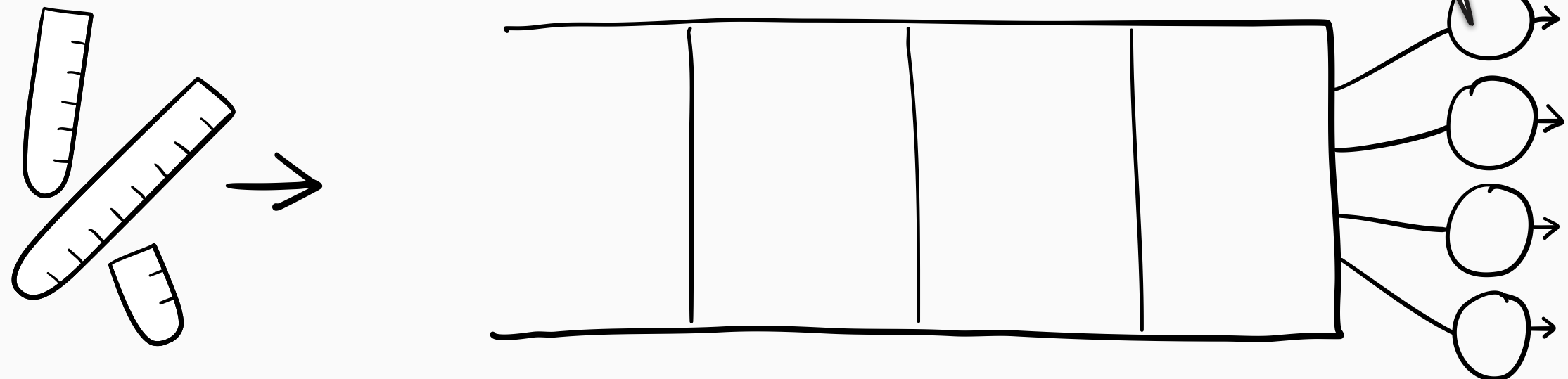
SRPT-1 (single-server): serves the job of least rank

similarly for Gittins-*k*, etc. PT: rank = remaining size

SRPT-*k* (multiserver): serves the *k* jobs of least rank

# Multiserver queueing system

M/G/ $k$



analyzed

SRPT-1 (**single-server**): serves the job of least **rank**

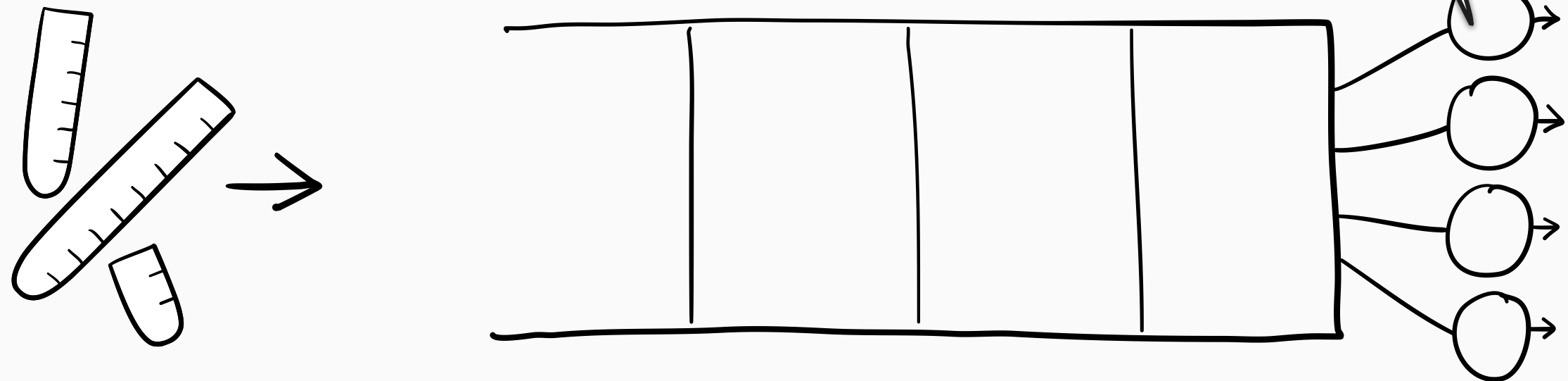
SRPT: **rank** = remaining size

SRPT- $k$  (**multiserver**): serves the  $k$  jobs of least **rank**



# Multiserver queueing system

M/G/k



analyzed

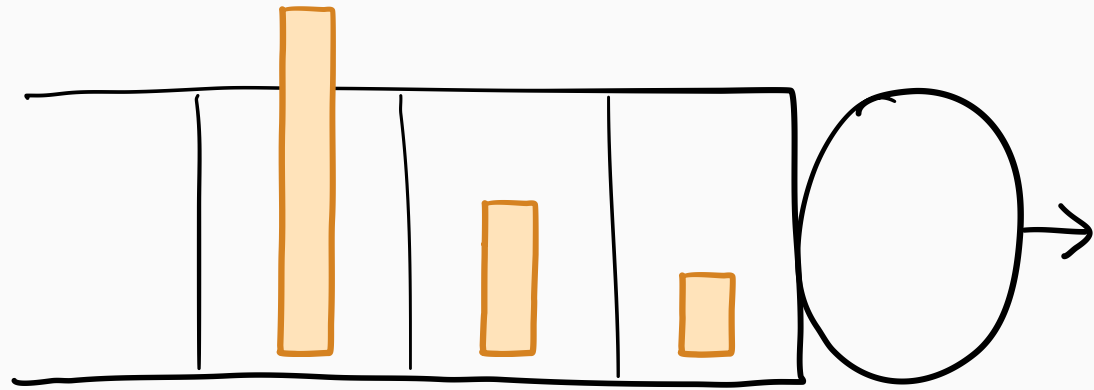
SRPT-1 (**single-server**): serves the job of least **rank**

! analysis open

**rank** = remaining size

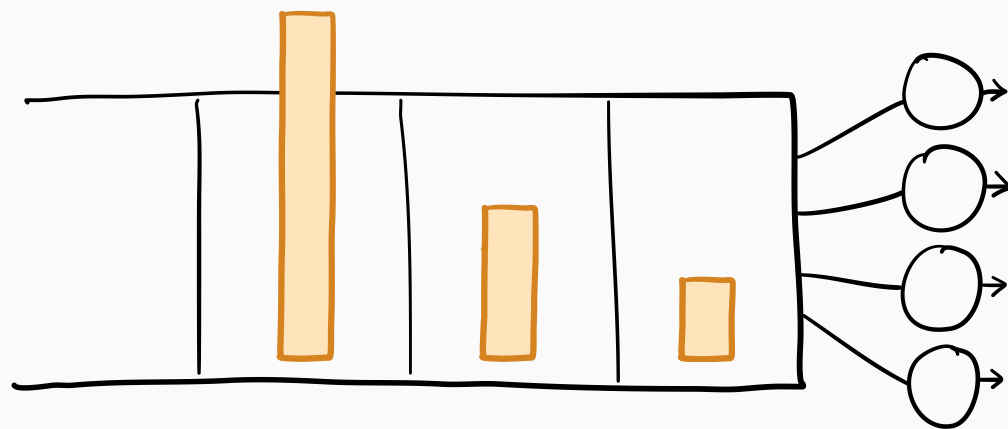
SRPT-**k** (**multiserver**): serves the **k** jobs of least **rank**

## Single-server system



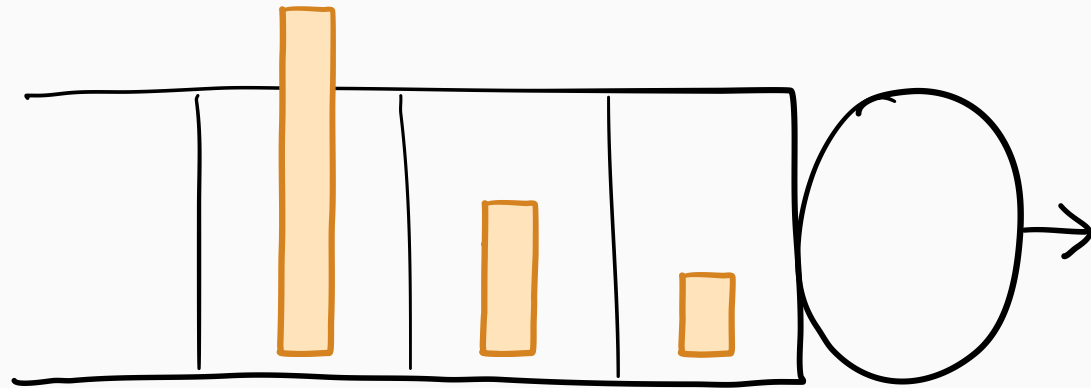
---

## Multiserver system



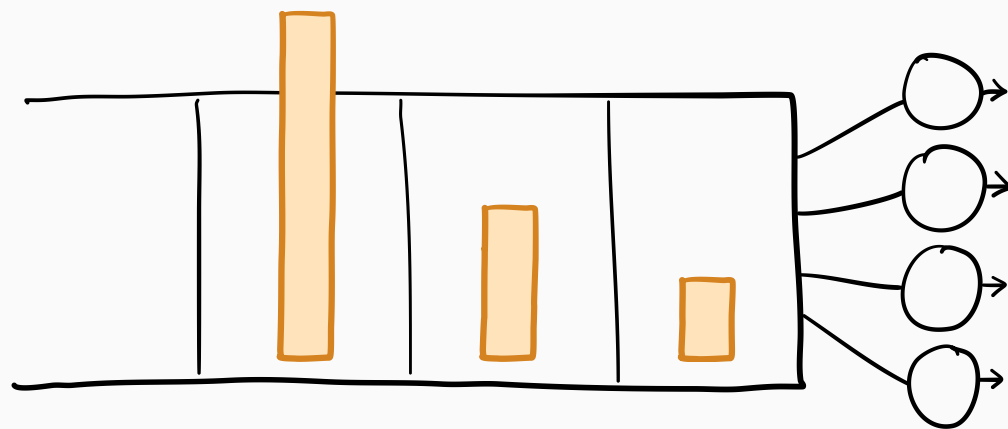
## Single-server system

server is “choke point”

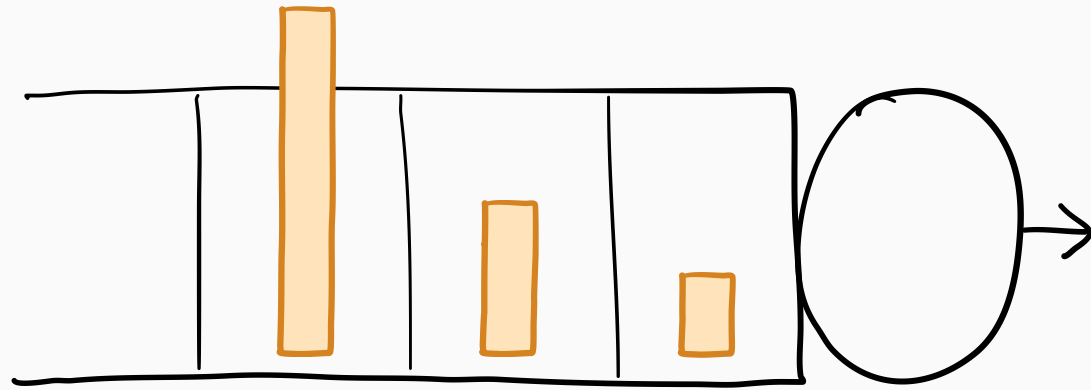


---

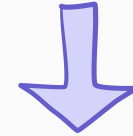
## Multiserver system



## Single-server system



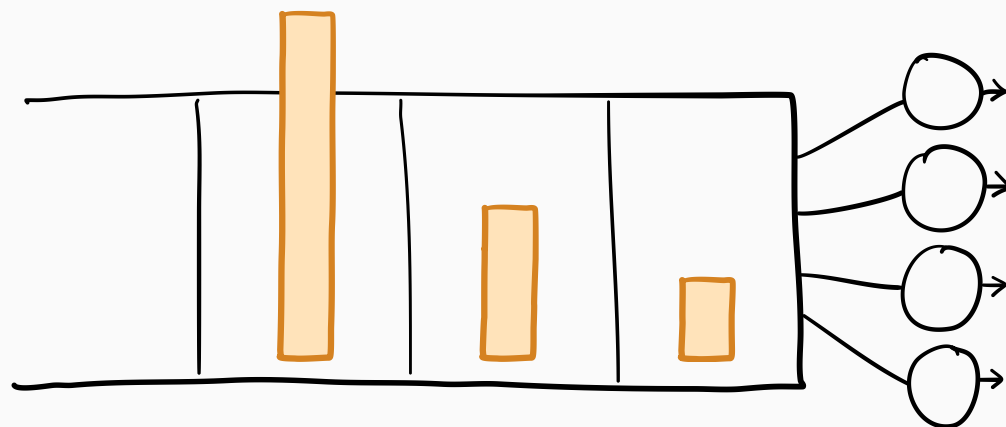
server is “choke point”



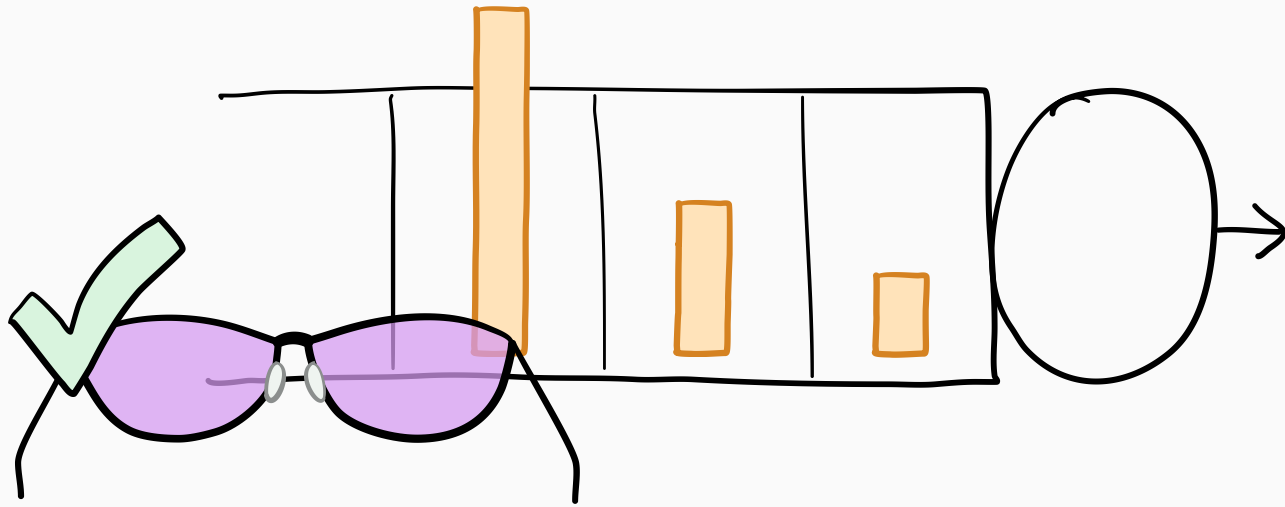
**rank** ordering absolute

---

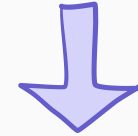
## Multiserver system



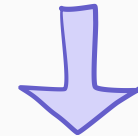
## Single-server system



server is “choke point”



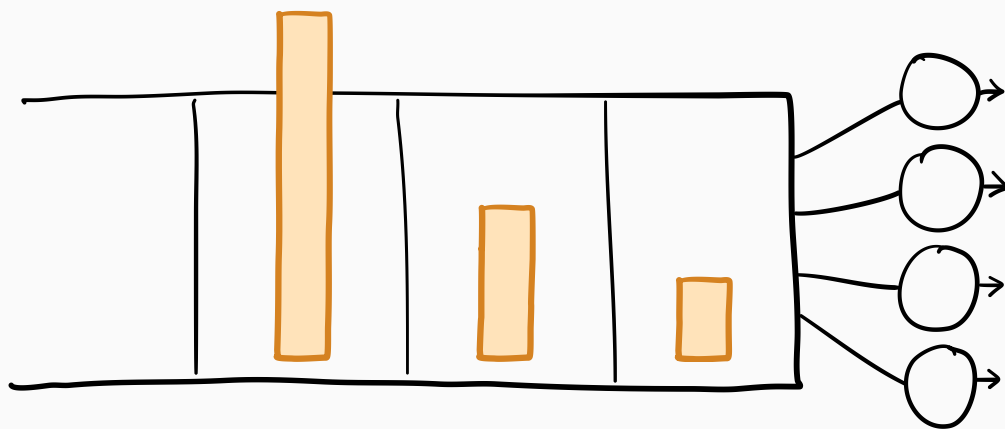
**rank** ordering absolute



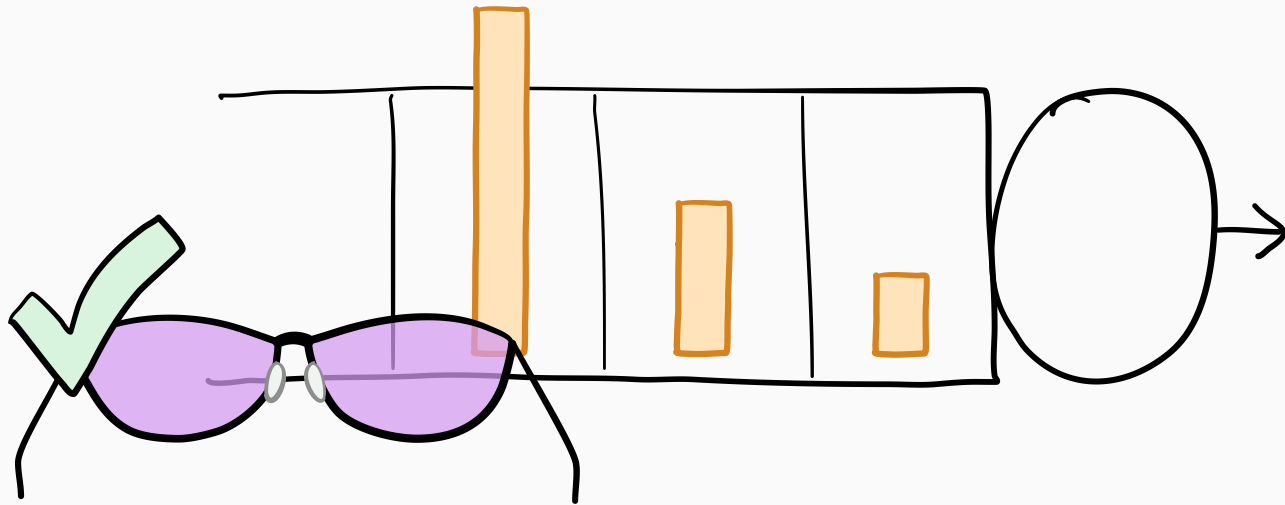
observed *r*-work determines *T*

---

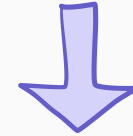
## Multiserver system



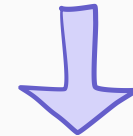
## Single-server system



server is “choke point”



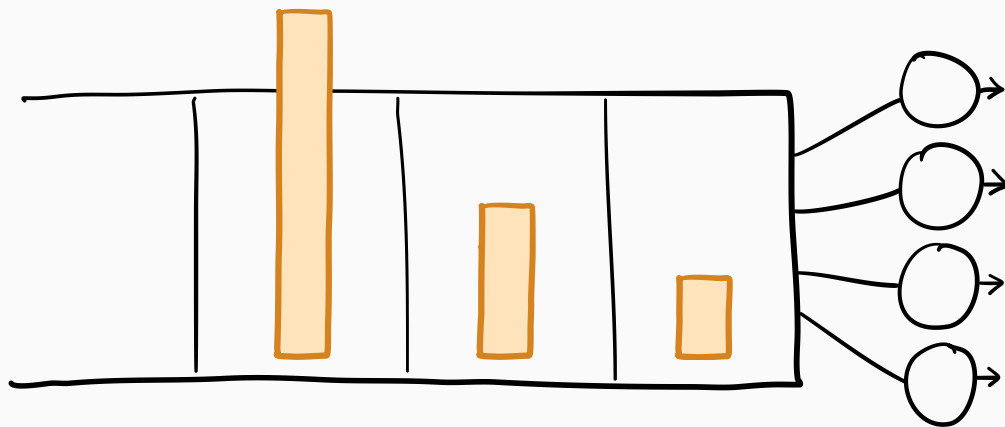
**rank** ordering absolute



observed *r*-work determines *T*

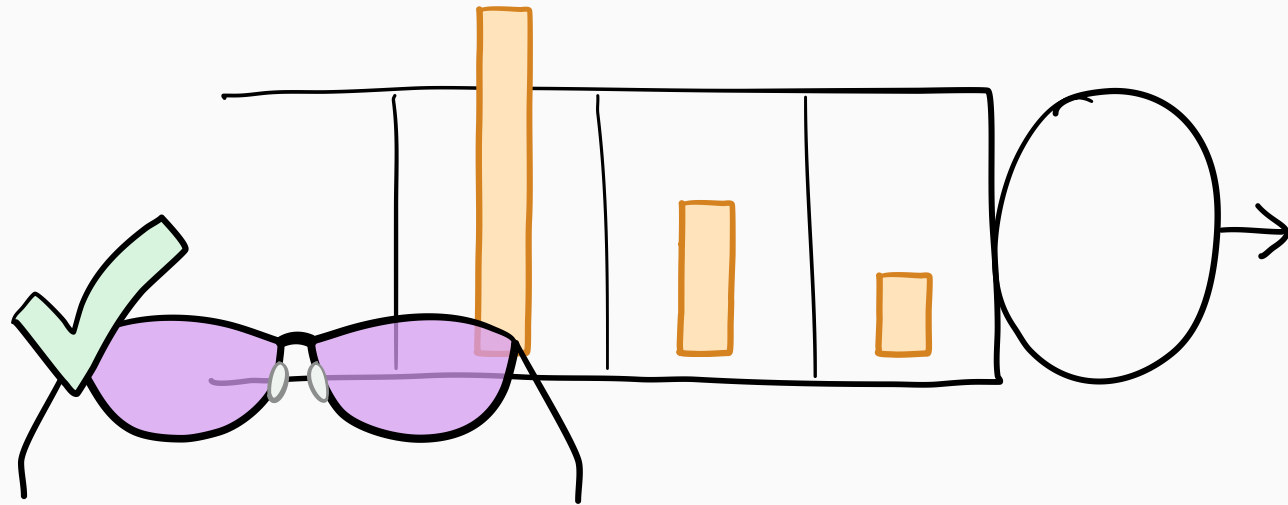
---

## Multiserver system

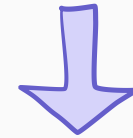


*no* single “choke point”

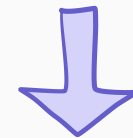
## Single-server system



server is “choke point”

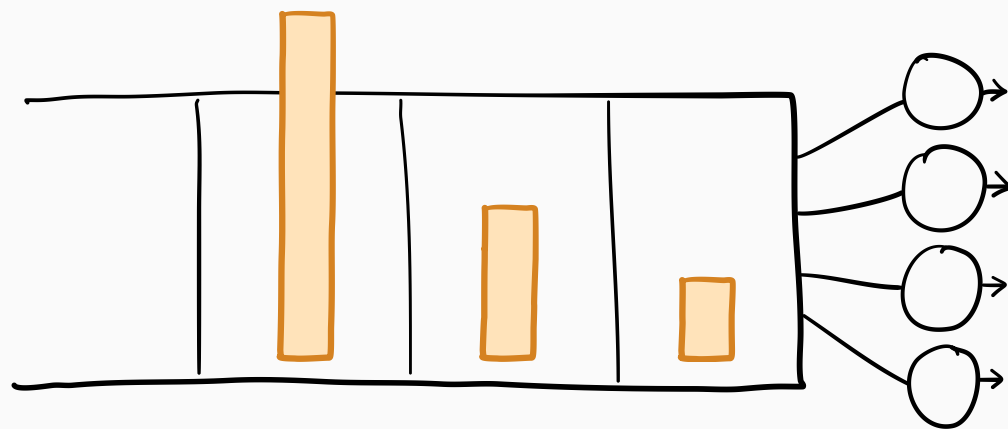


**rank** ordering absolute

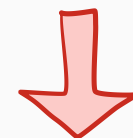


observed  $r$ -work determines  $T$

## Multiserver system

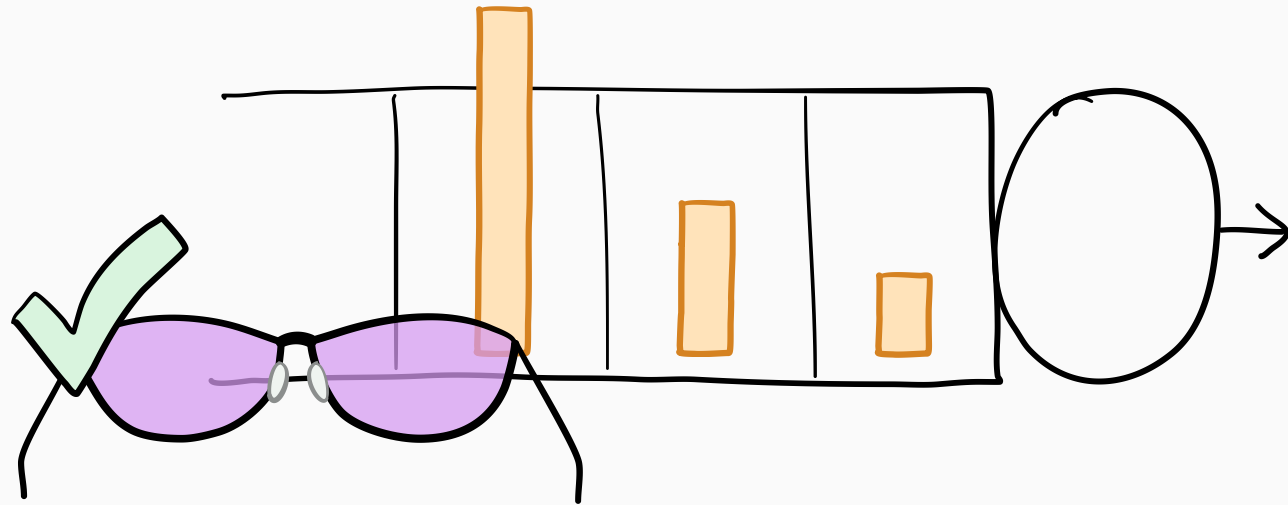


*no* single “choke point”

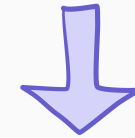


**rank** ordering *not* absolute

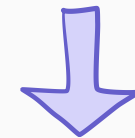
## Single-server system



server is “choke point”

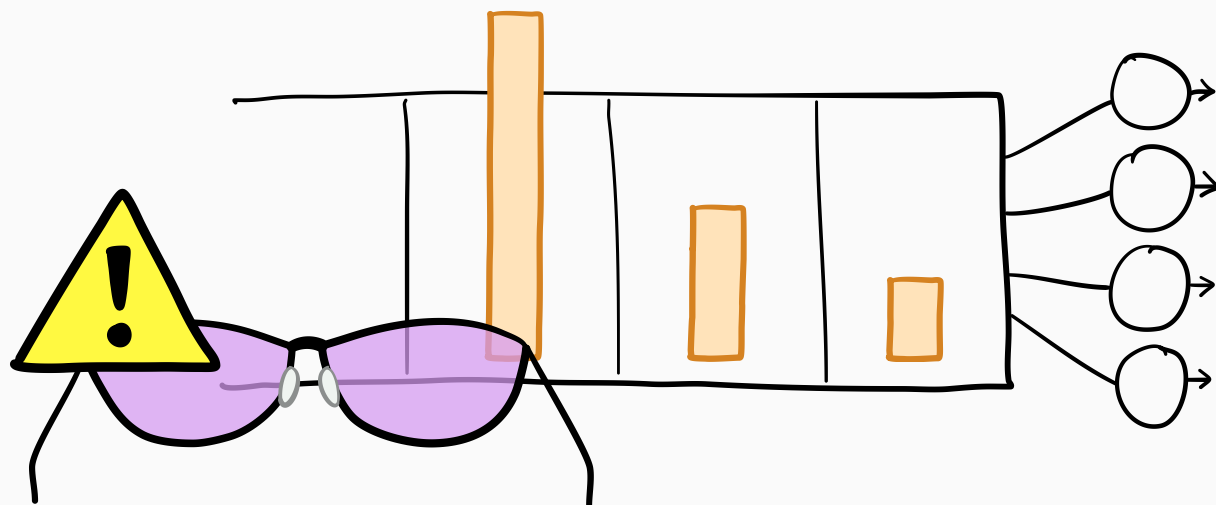


**rank** ordering absolute

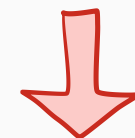


observed  $r$ -work determines  $T$

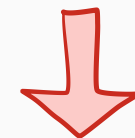
## Multiserver system



*no* single “choke point”

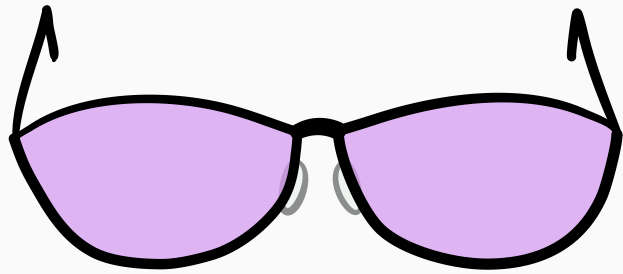


**rank** ordering *not* absolute

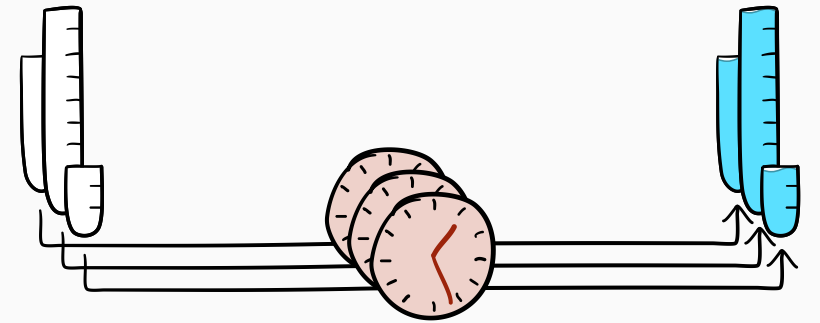


observed  $r$ -work *not enough!*

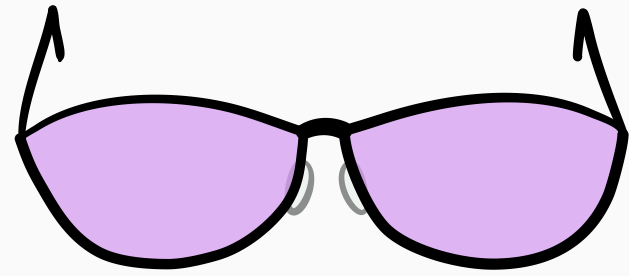




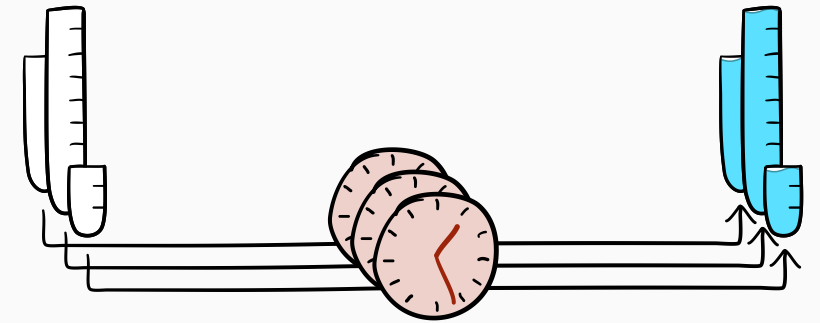
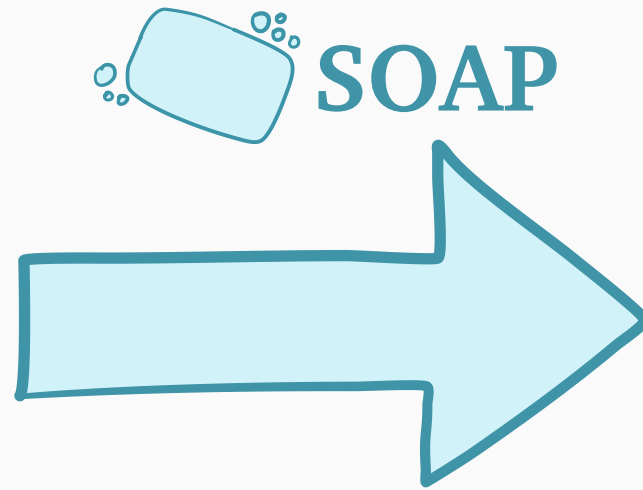
*r*-work  $W(r)$



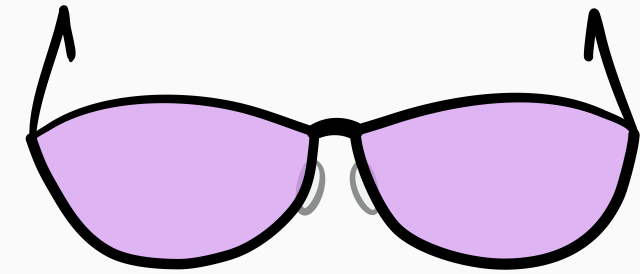
response time  $T$



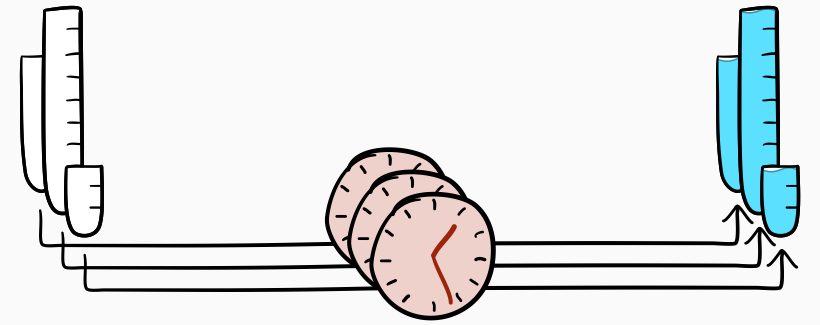
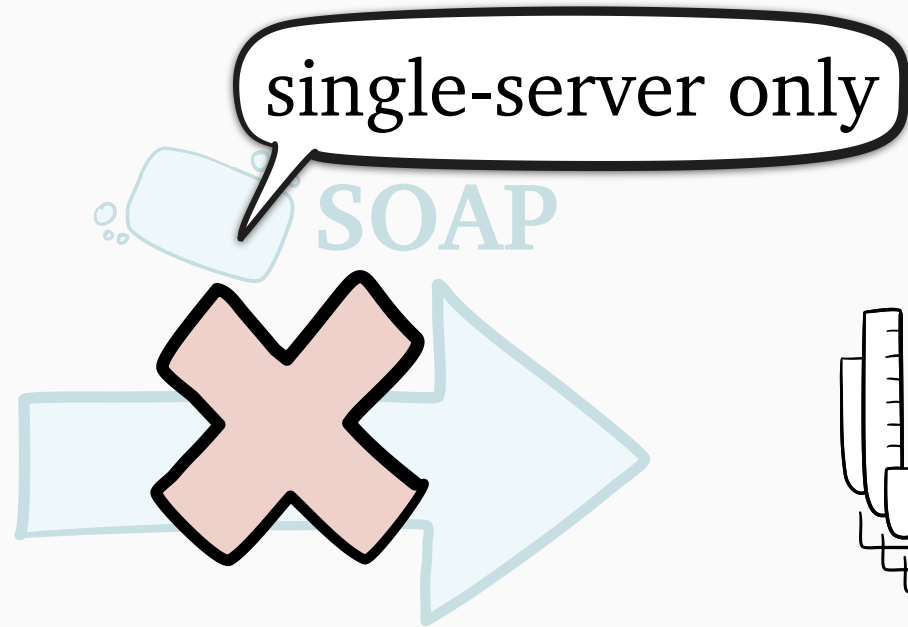
*r*-work  $W(r)$



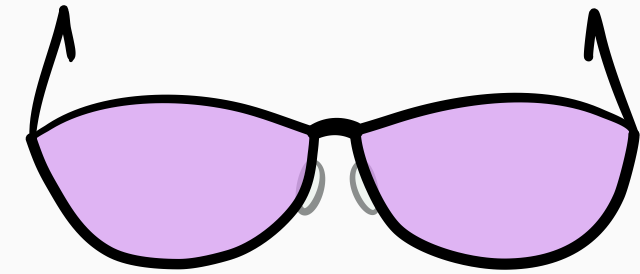
response time  $T$



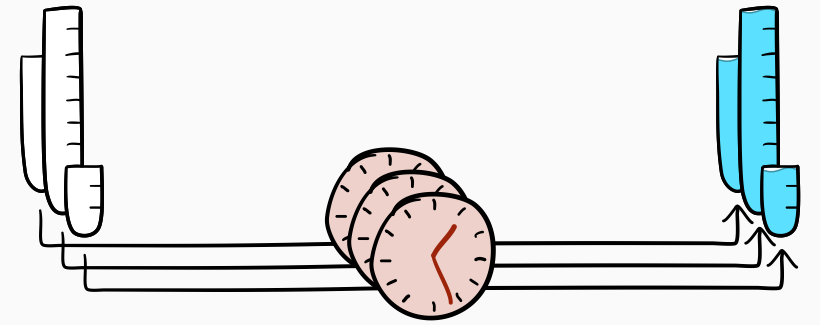
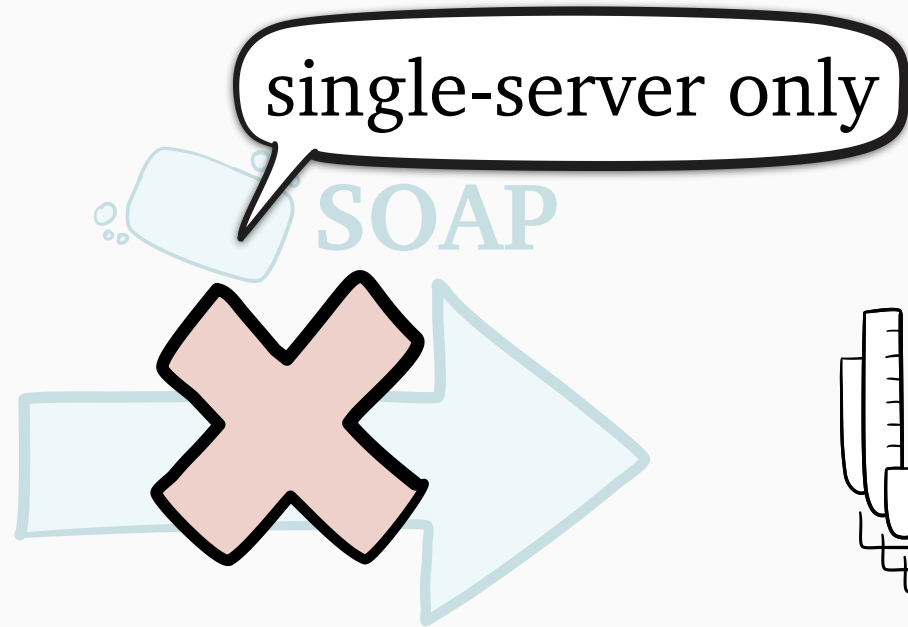
*r*-work  $W(r)$



response time  $T$



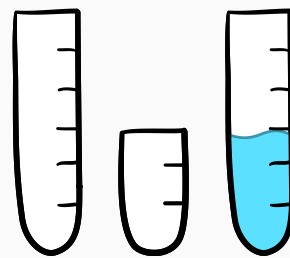
$r$ -work  $W(r)$



response time  $T$



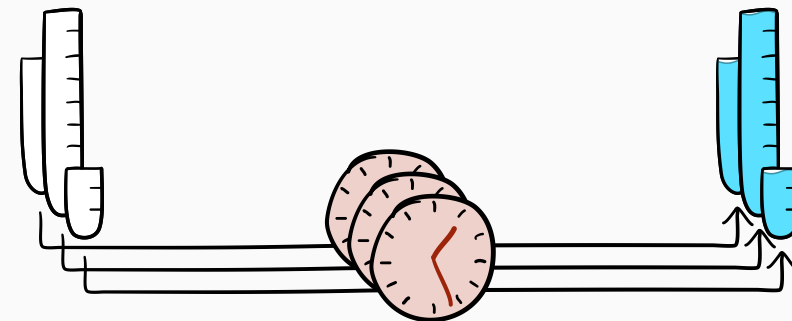
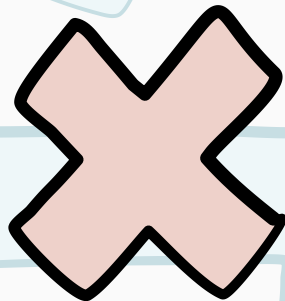
Little's law



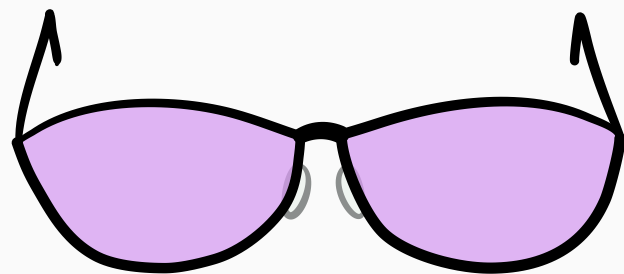
number of jobs  $N$

single-server only

SOAP



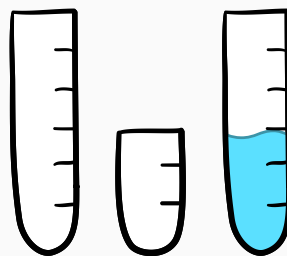
response time  $T$



$r$ -work  $W(r)$

$$E[N] = \lambda E[T]$$

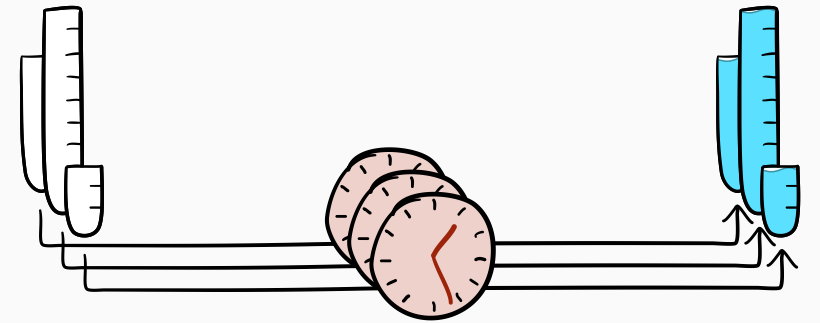
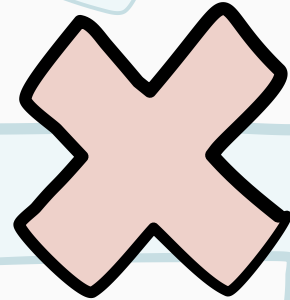
Little's law



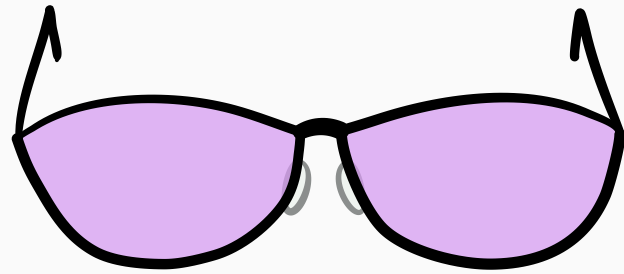
number of jobs  $N$

single-server only

SOAP



response time  $T$

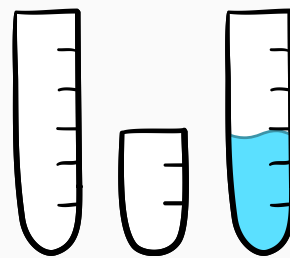


$r$ -work  $W(r)$

$$E[N] = \lambda E[T]$$

Little's law

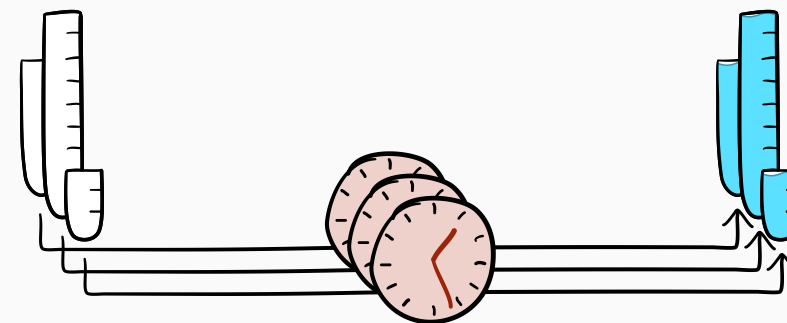
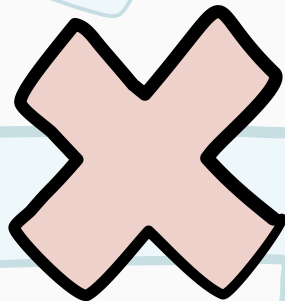
any number of servers



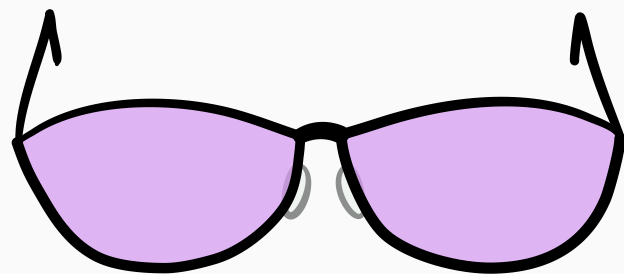
number of jobs  $N$

single-server only

SOAP



response time  $T$



$r$ -work  $W(r)$



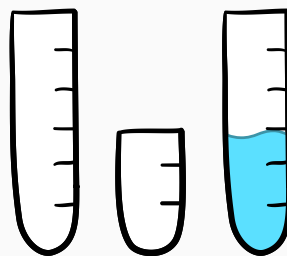
WINE



$$E[N] = \lambda E[T]$$

Little's law

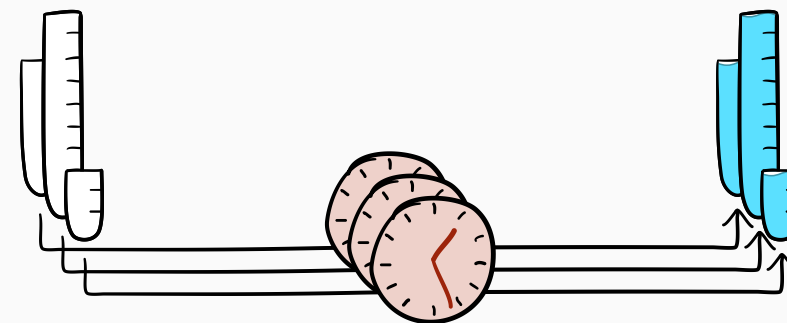
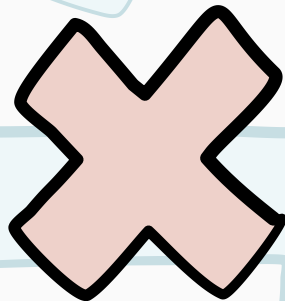
any number of servers



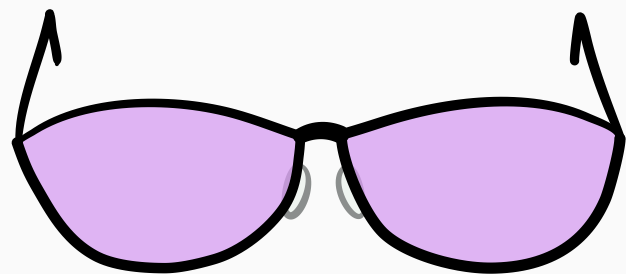
number of jobs  $N$

single-server only

SOAP



response time  $T$



$r$ -work  $W(r)$



WINE

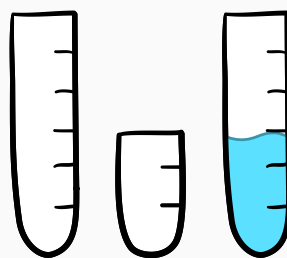
any number  
of servers



$$E[N] = \lambda E[T]$$

Little's law

any number  
of servers



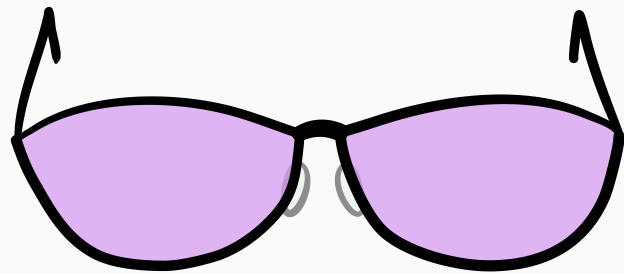
number of jobs  $N$



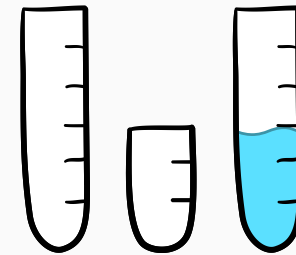
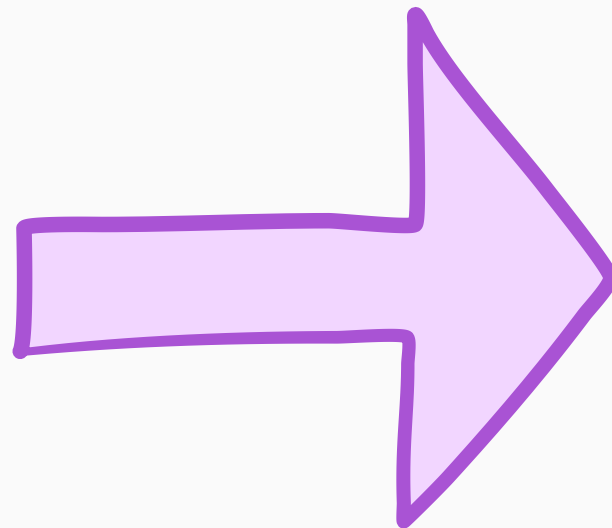


# WINE

Work Integral Number Equality



*r*-work  $W(r)$

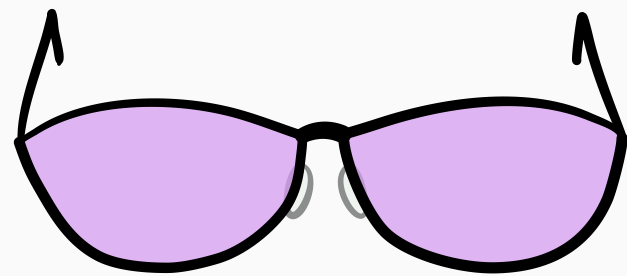


number of jobs  $N$

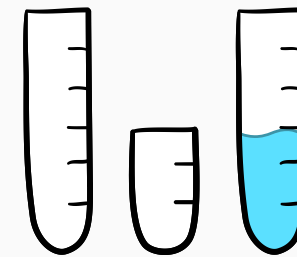
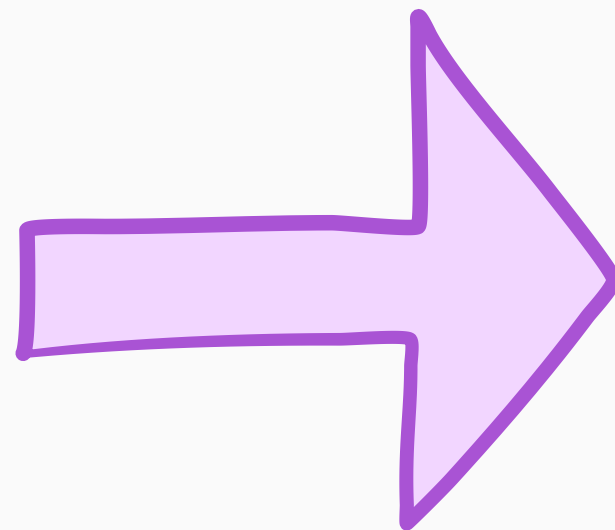


# WINE

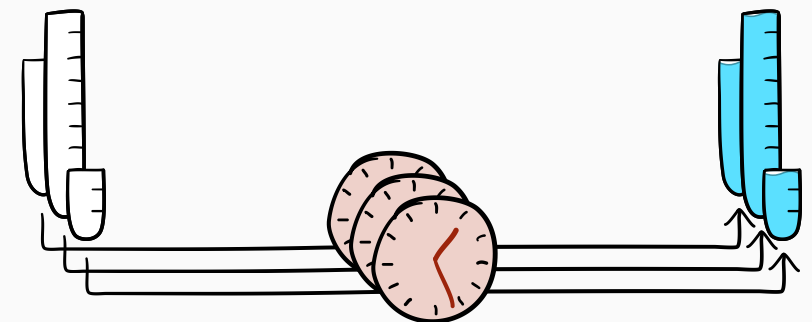
Work Integral Number Equality



$r$ -work  $W(r)$



number of jobs  $N$

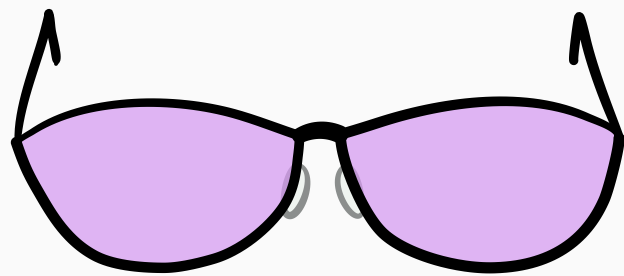


mean response time  $E[T]$



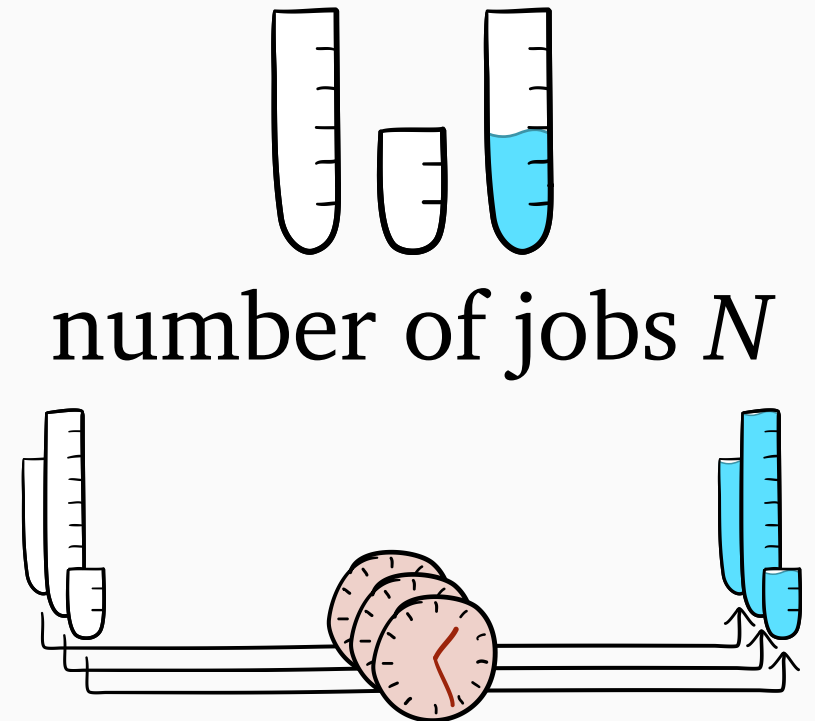
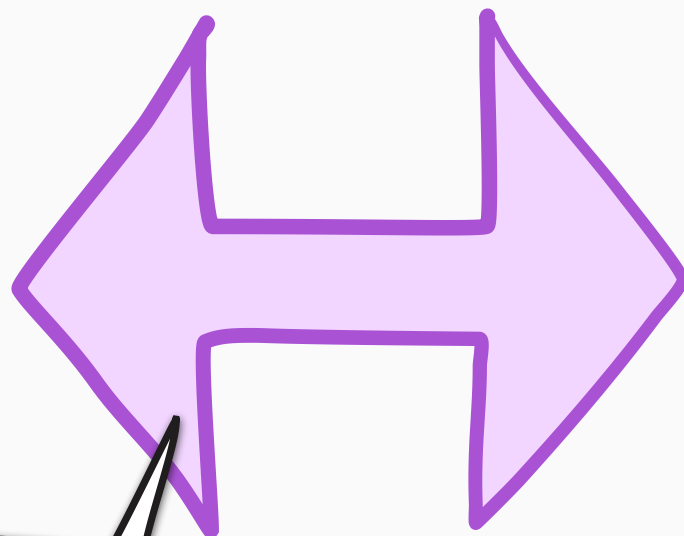
# WINE

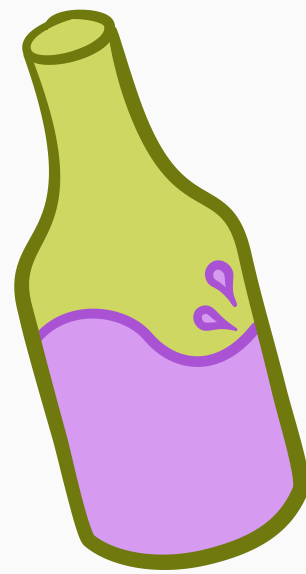
Work Integral Number Equality



$r$ -work  $W(r)$

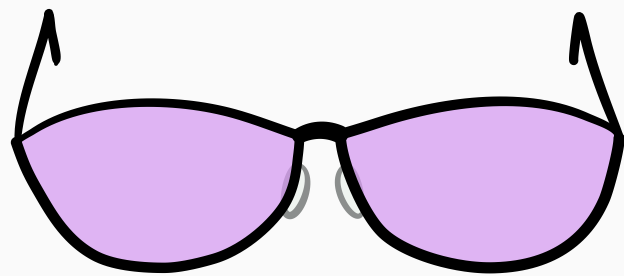
reduces  $E[T]$   
to  $r$ -work



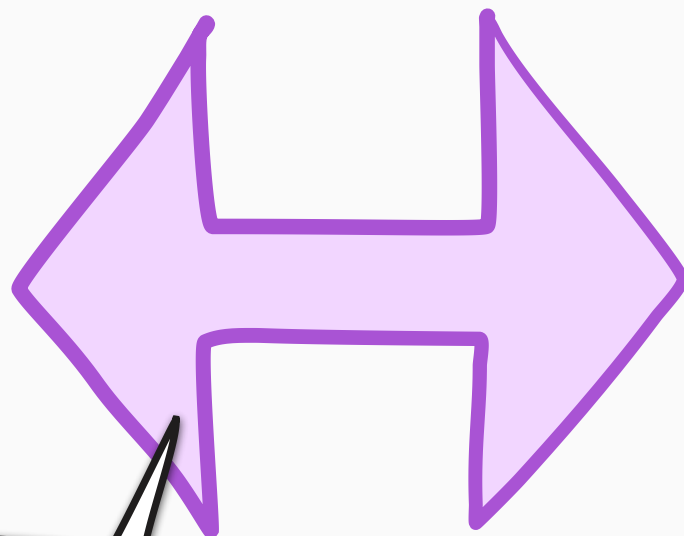


# WINE

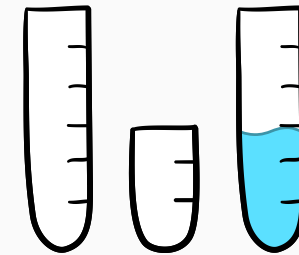
Work Integral Number Equality



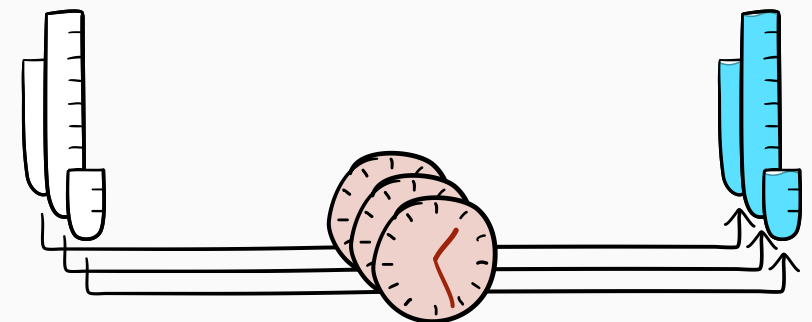
$r$ -work  $W(r)$



reduces  $E[T]$   
to  $r$ -work



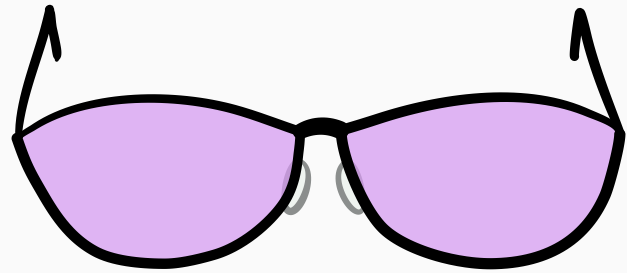
number of jobs  $N$



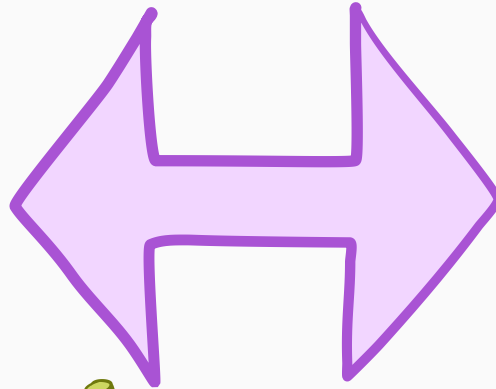
mean response time  $E[T]$



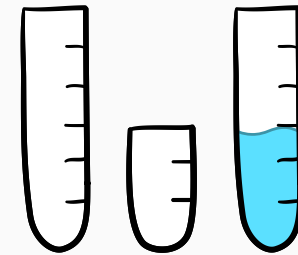
**NEW!** First analysis of SRPT- $k$ , Gittins- $k$ , noisy size estimates



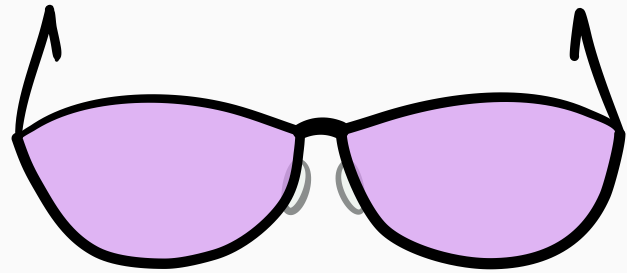
*r*-work  $W(r)$



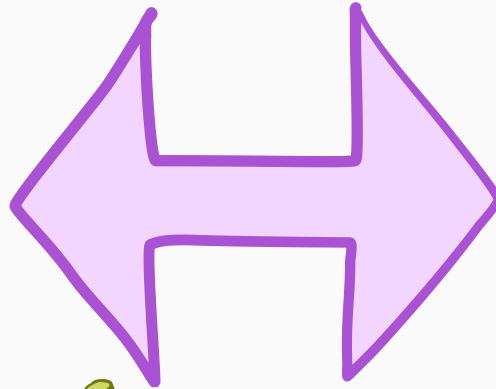
WINE



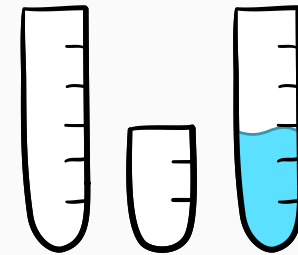
number of jobs  $N$



*r*-work  $W(r)$



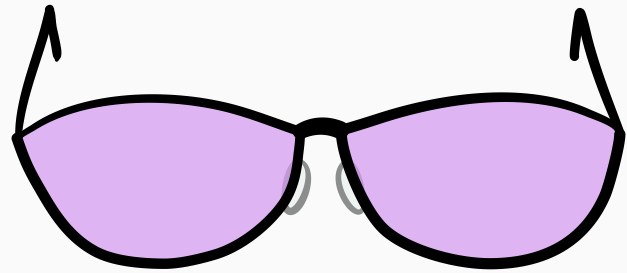
WINE



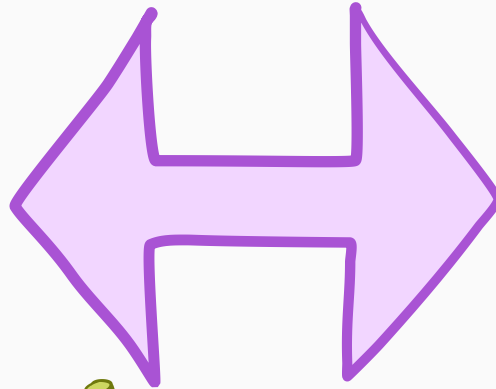
number of jobs  $N$



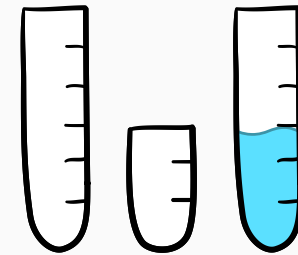
What is *r*-work?



*r*-work  $W(r)$



WINE



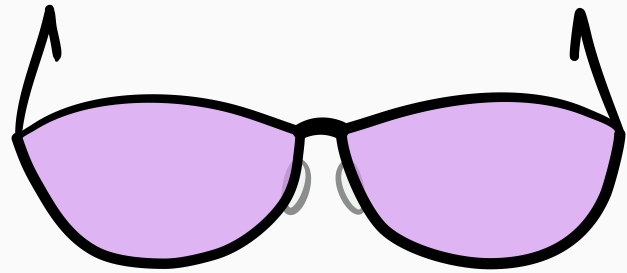
number of jobs  $N$



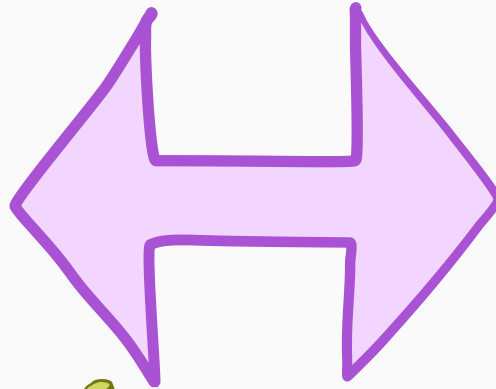
What is *r*-work?



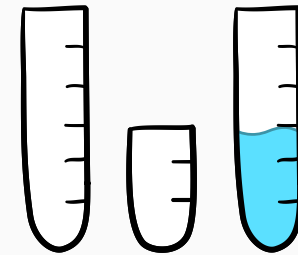
How do we get number of jobs from *r*-work?



*r*-work  $W(r)$



WINE



number of jobs  $N$



What is *r*-work?

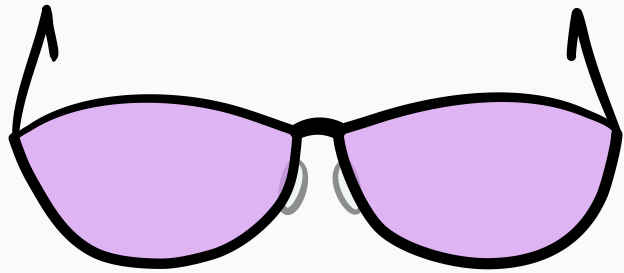


How do we get number of jobs from *r*-work?

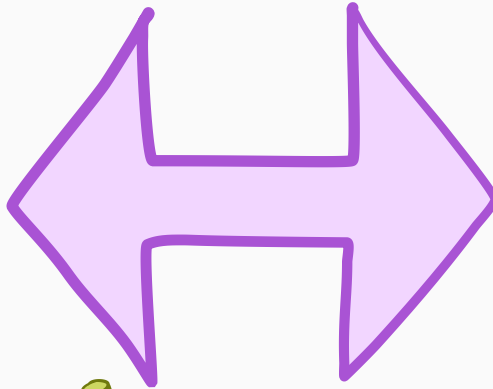


How do we analyze *r*-work?

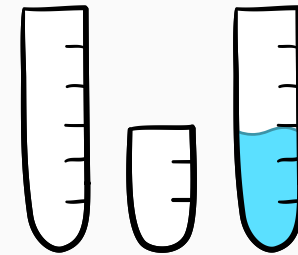




*r*-work  $W(r)$



WINE



number of jobs  $N$

This talk:  
SRPT-*k*



What is *r*-work?



How do we get number of jobs from *r*-work?



How do we analyze *r*-work?

# Defining $r$ -work

for SRPT

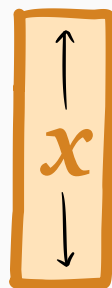
$W(r)$  = work relevant to **rank**  $r$

# Defining $r$ -work

for SRPT

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

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

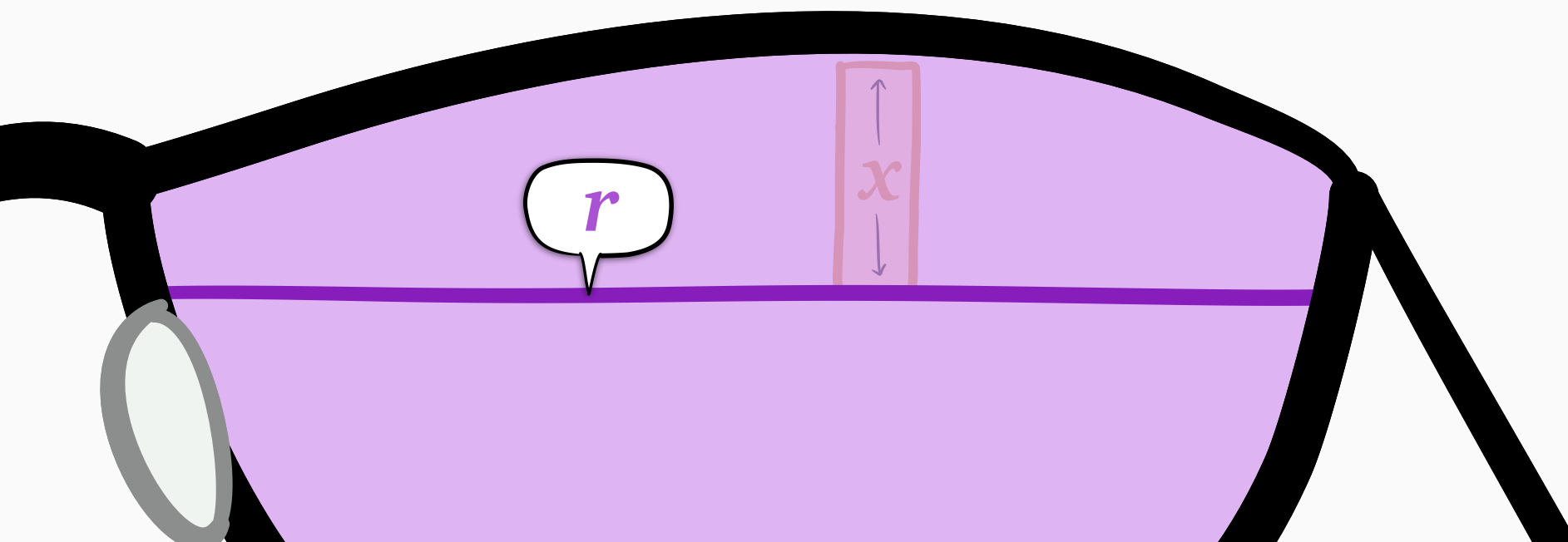


# Defining $r$ -work

for SRPT

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

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

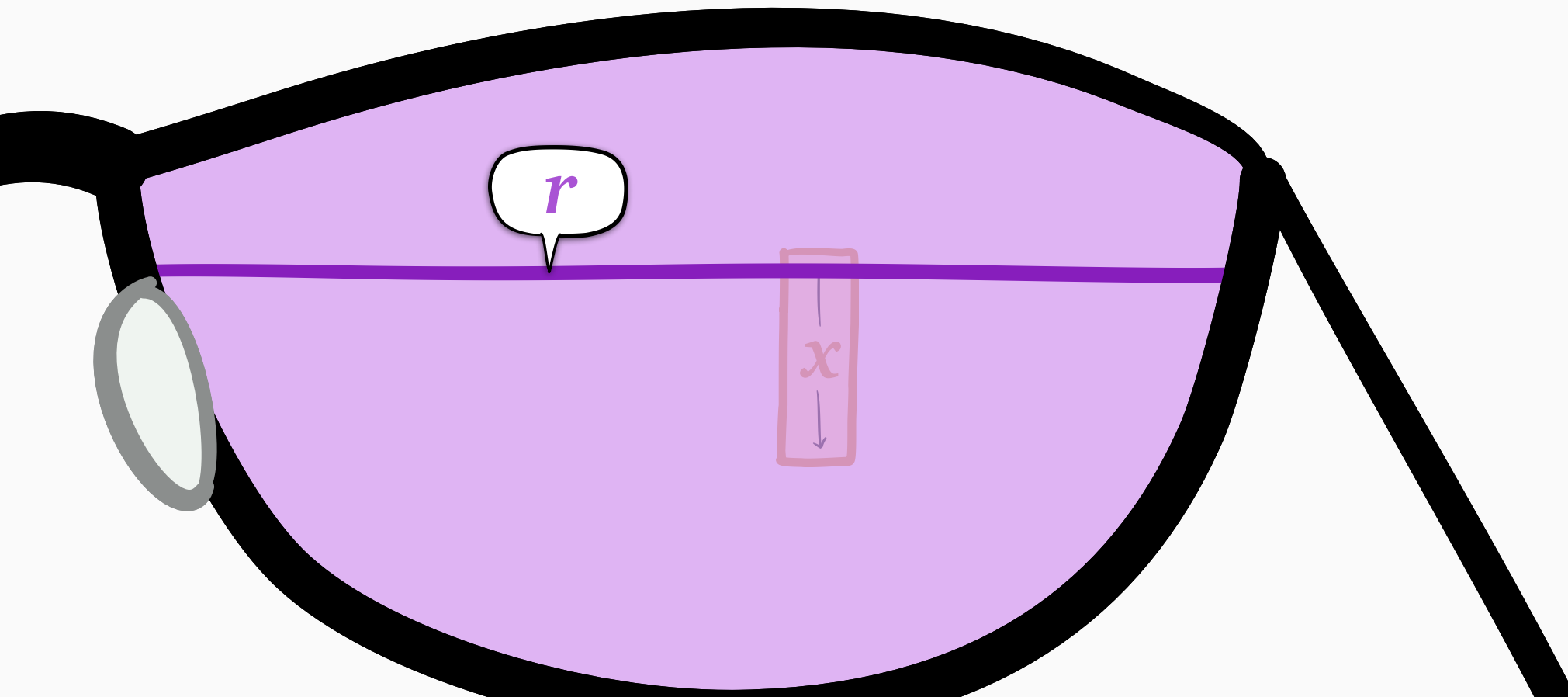


# Defining $r$ -work

for SRPT

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

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

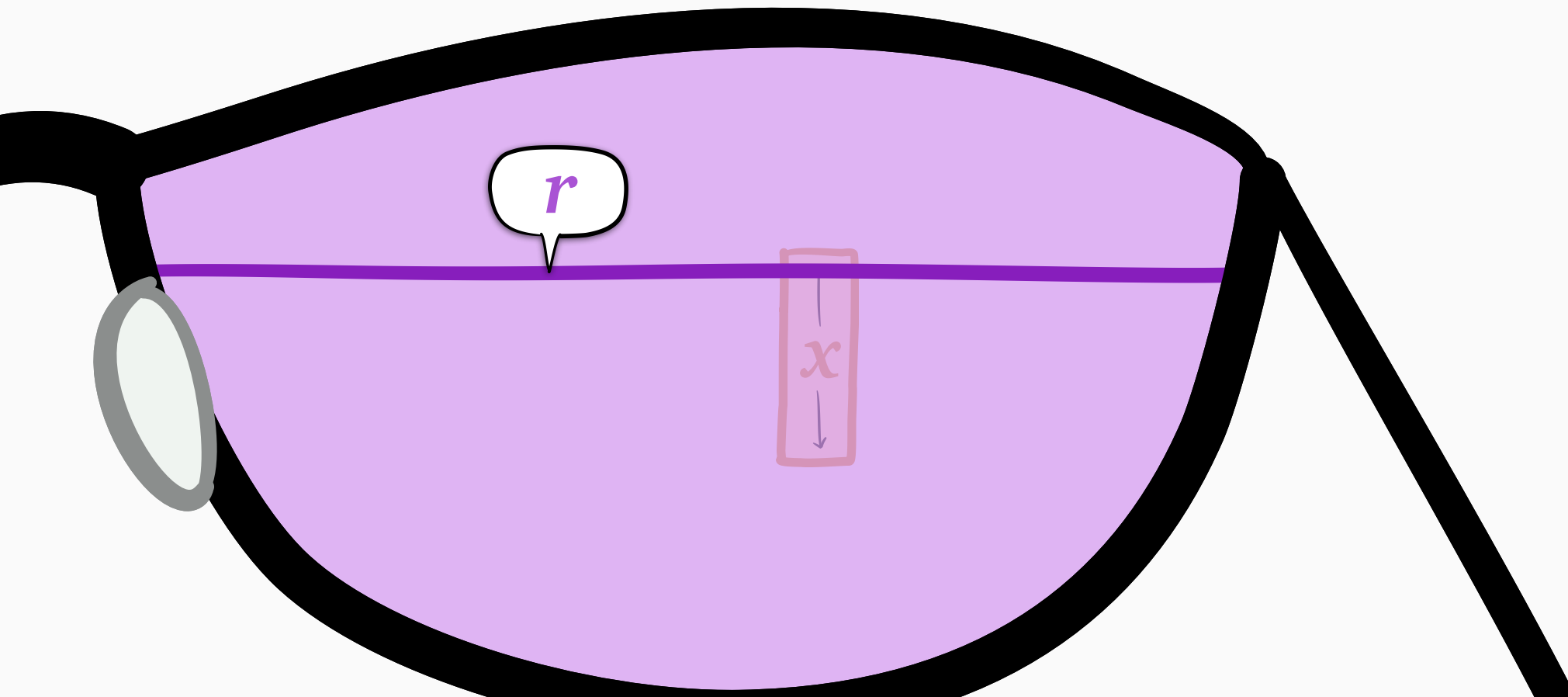


# Defining $r$ -work

for SRPT

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

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

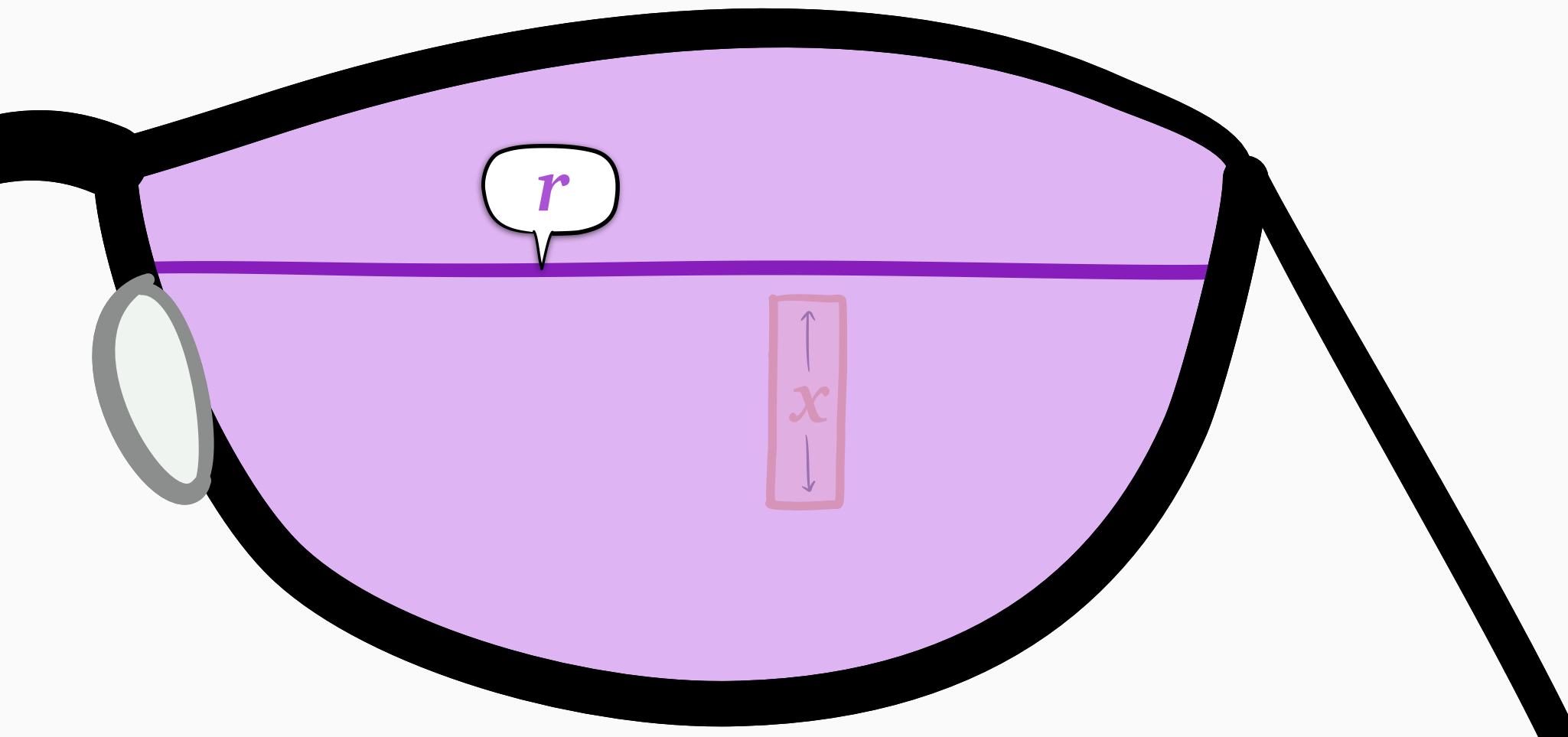


# Defining $r$ -work

for SRPT

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

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

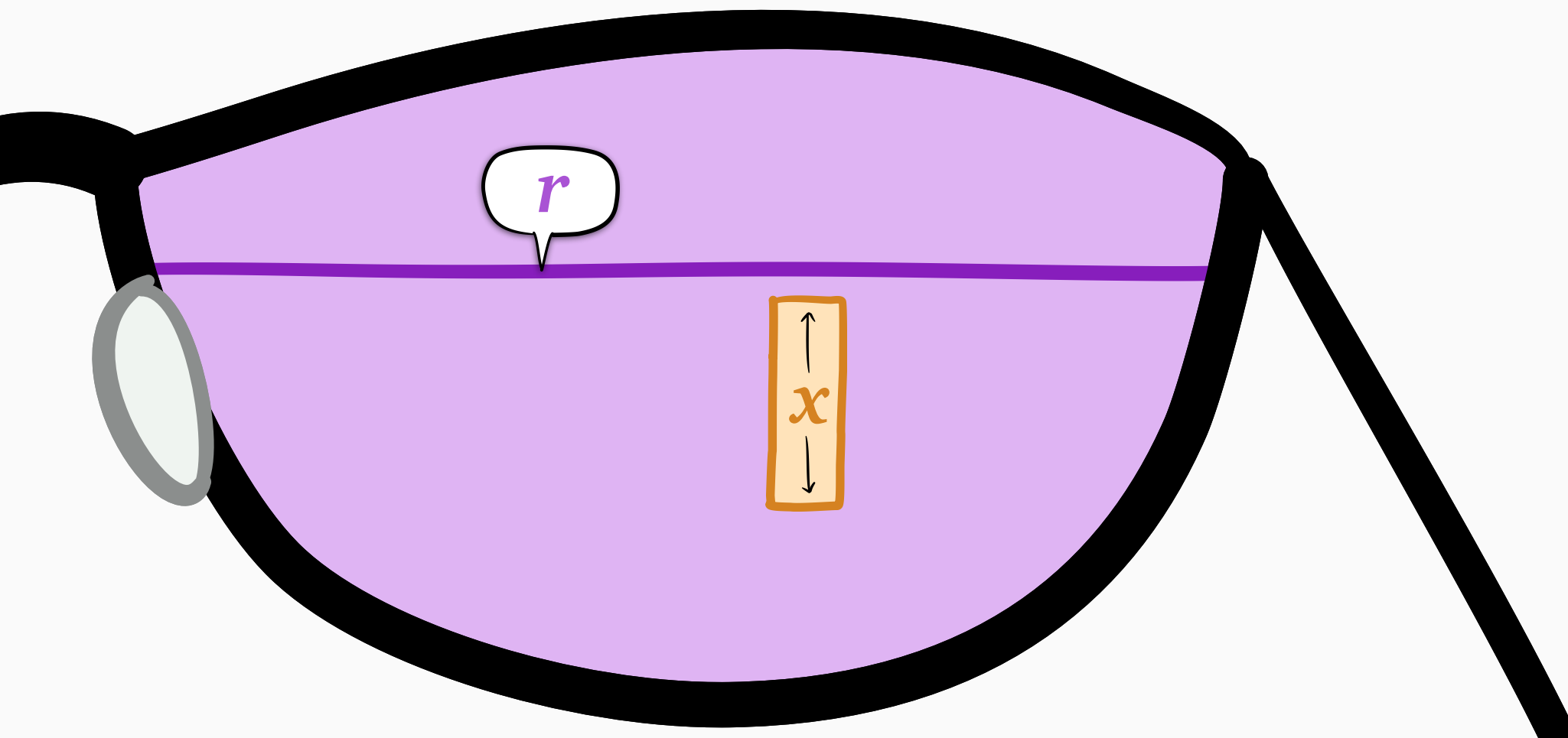


# Defining $r$ -work

for SRPT

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

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



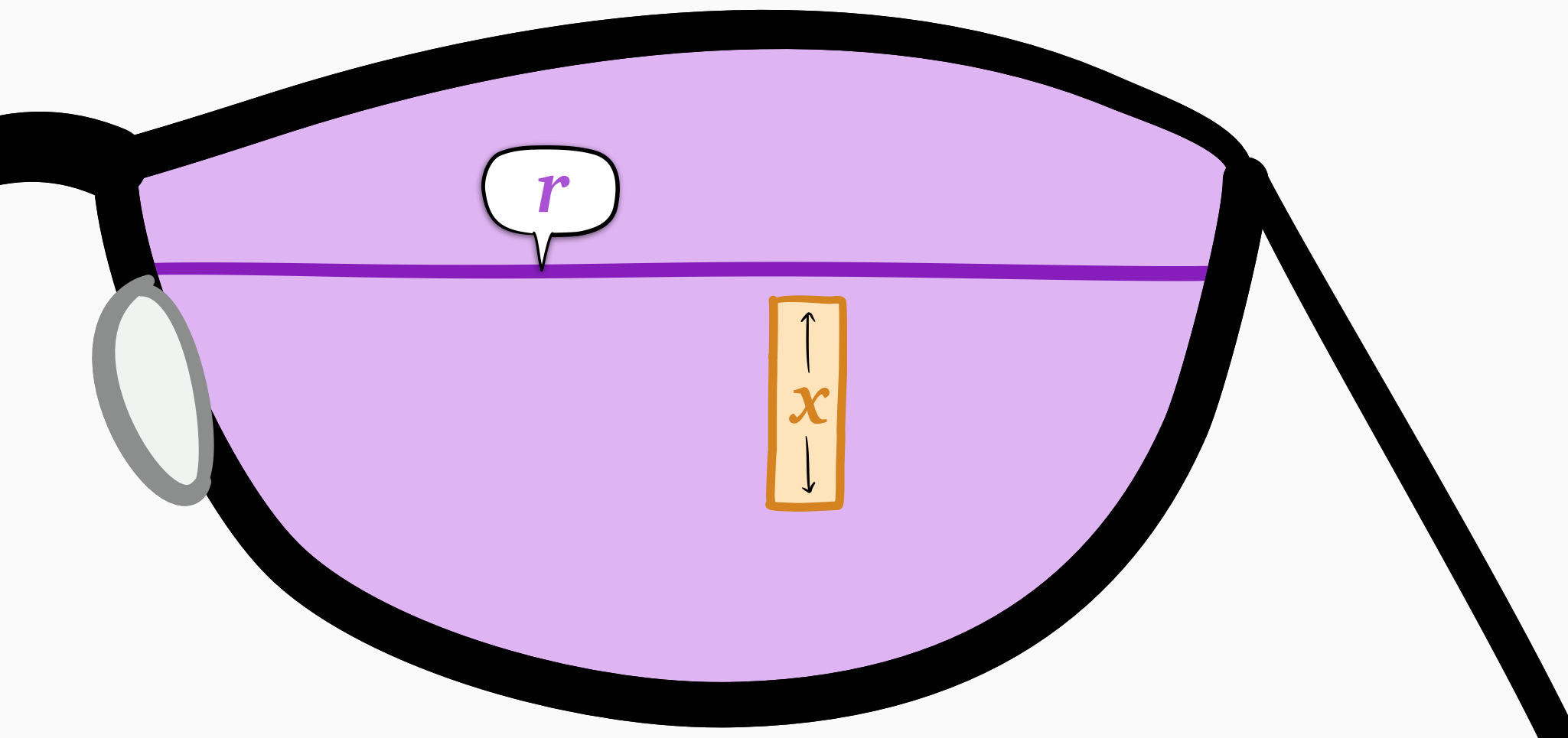


# Defining $r$ -work

for SRPT

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

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

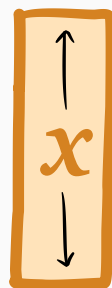


# Defining $r$ -work

for SRPT

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

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

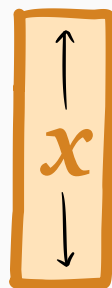


# Defining $r$ -work

for SRPT

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

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



From  $r$ -work to number of jobs  $N$

# From $r$ -work to number of jobs $N$

Goal: integral =  $N$

$W(r)$



# From $r$ -work to number of jobs $N$

**Goal:** integral =  $N$

$W(r)$



**Suffices:** integral = 1

$w_x(r)$

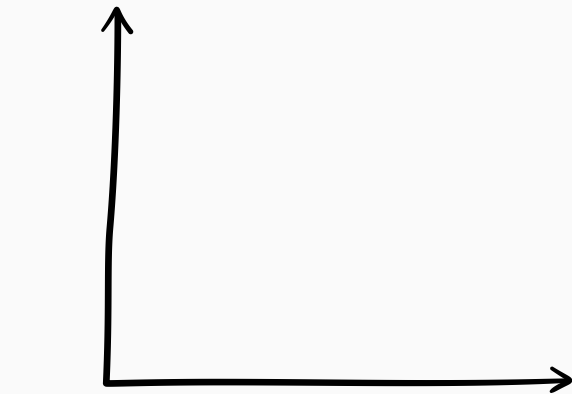


$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

# From $r$ -work to number of jobs $N$

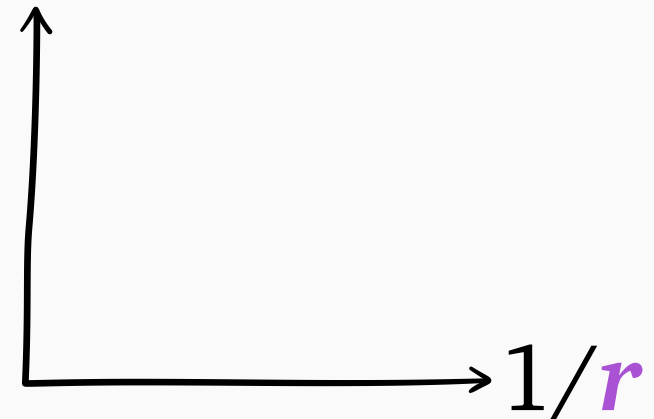
**Goal:** integral =  $N$

$W(r)$



**Suffices:** integral = 1

$w_x(r)$



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

# From $r$ -work to number of jobs $N$

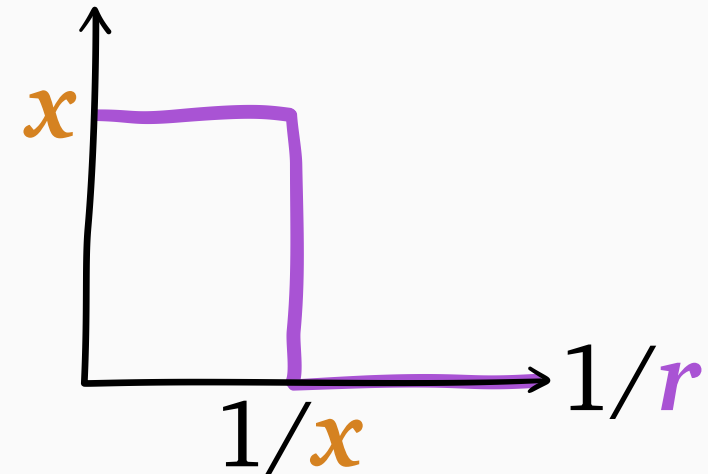
Goal: integral =  $N$

$W(r)$



Suffices: integral = 1

$w_x(r)$



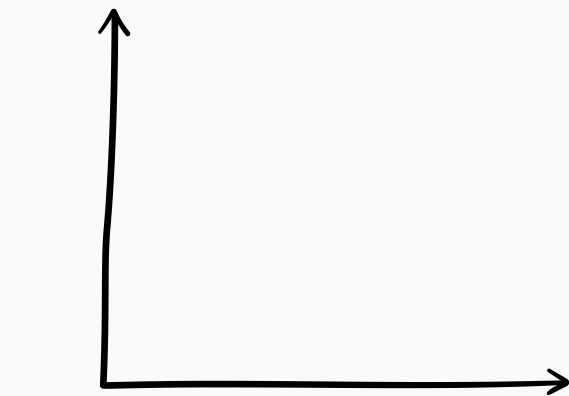
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$



# From $r$ -work to number of jobs $N$

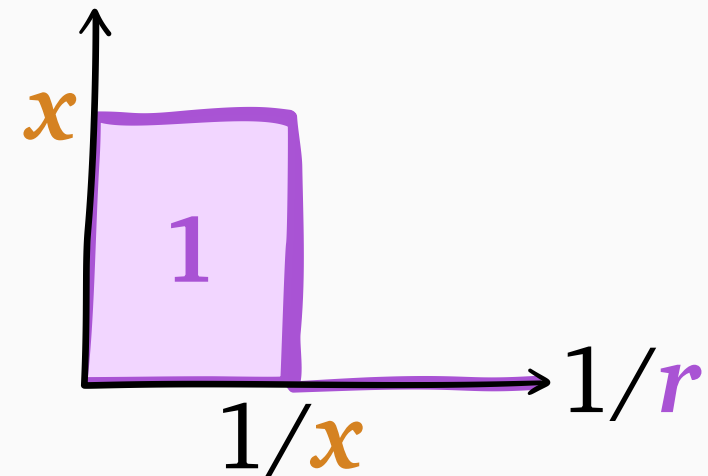
Goal: integral =  $N$

$W(r)$



Suffices: integral = 1

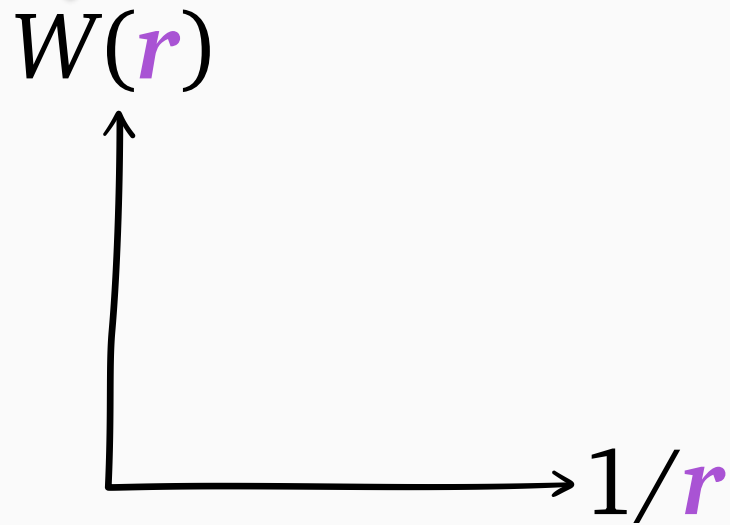
$w_x(r)$



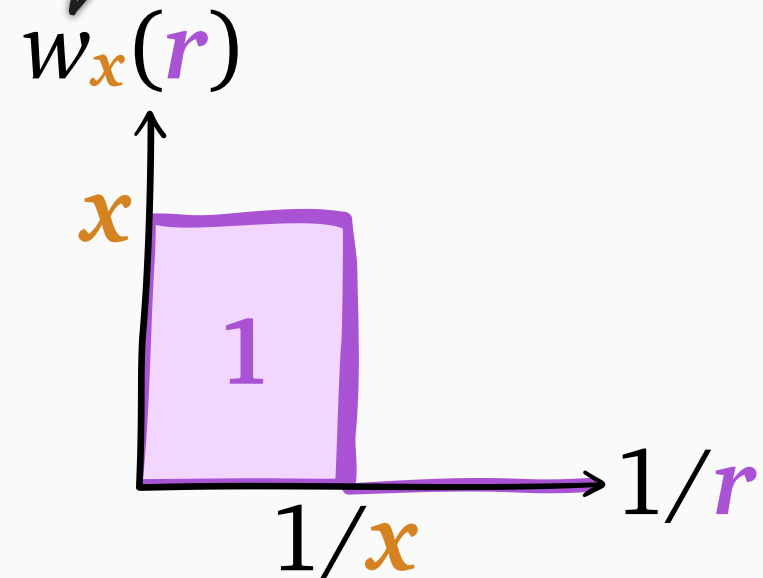
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

# From $r$ -work to number of jobs $N$

Goal: integral =  $N$



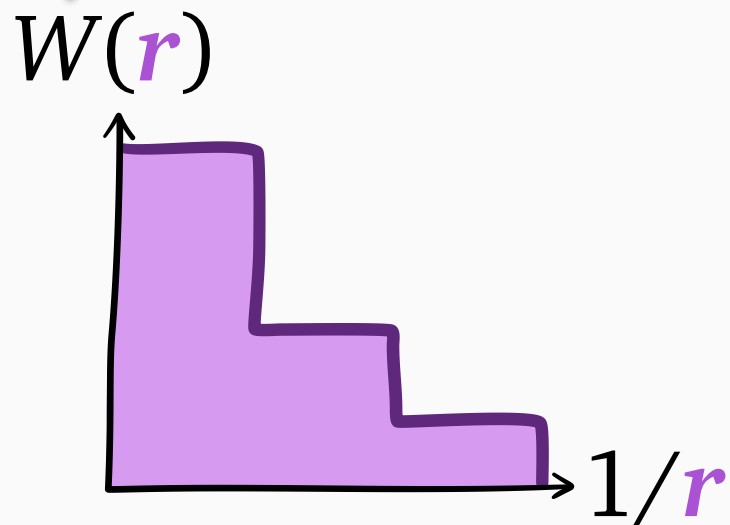
Suffices: integral = 1



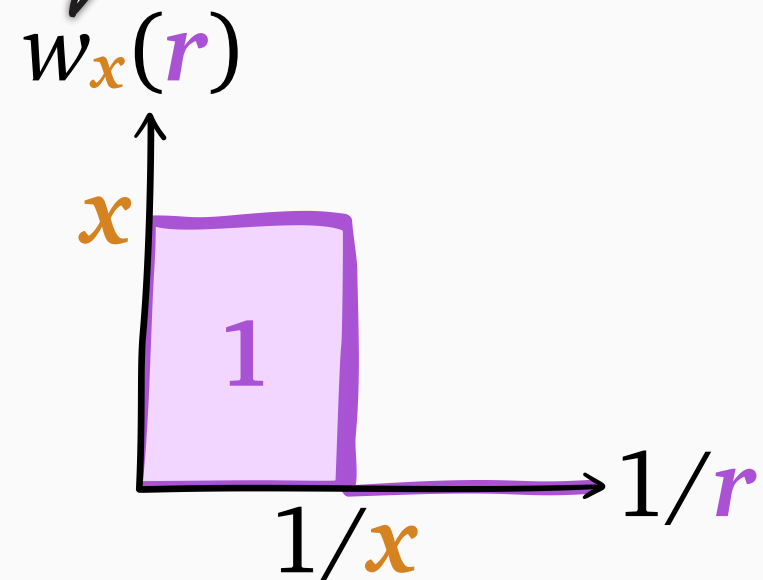
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

# From $r$ -work to number of jobs $N$

Goal: integral =  $N$



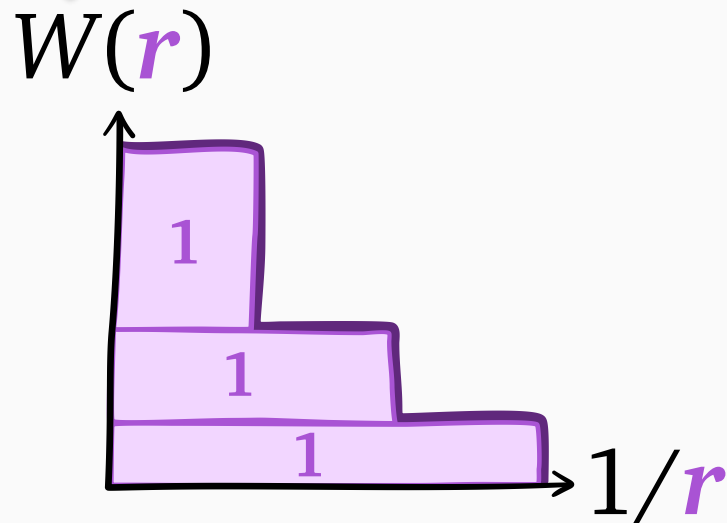
Suffices: integral = 1



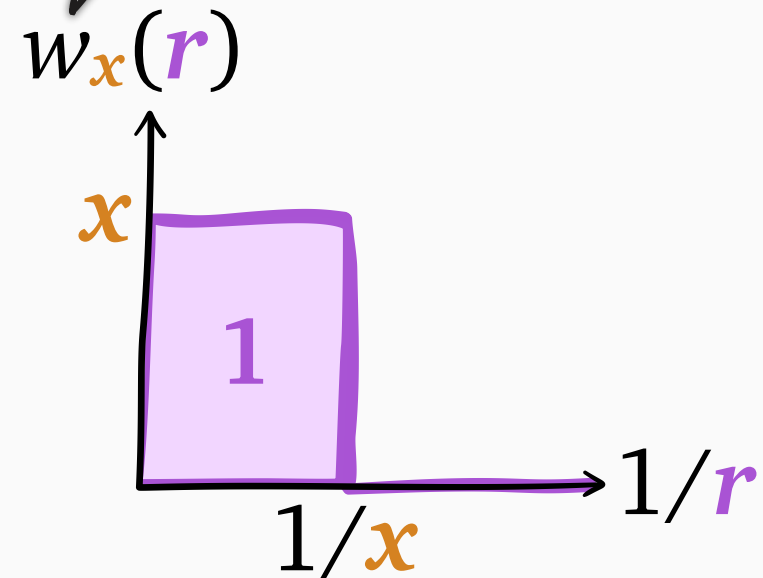
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

# From $r$ -work to number of jobs $N$

**Goal:** integral =  $N$



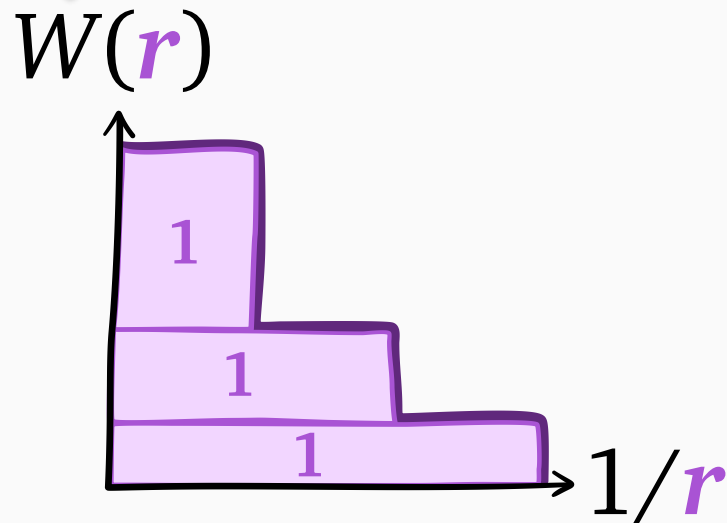
**Suffices:** integral = 1



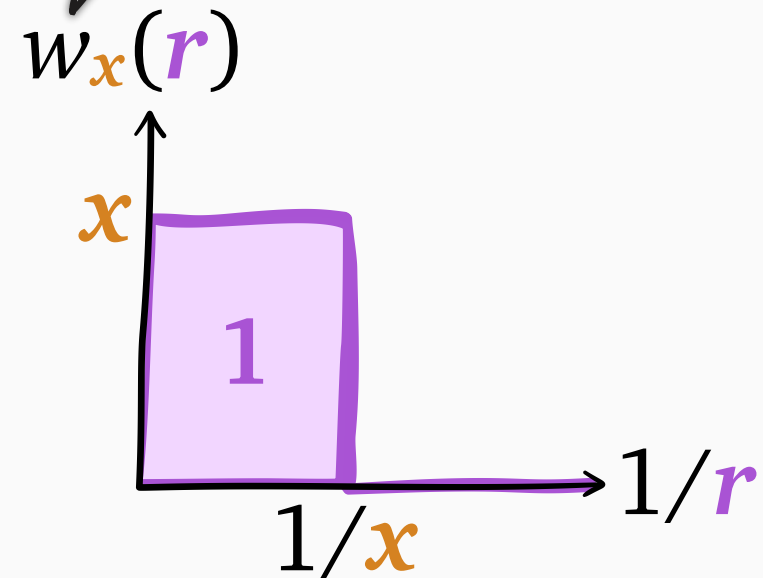
$$w_x(r) = r\text{-work of job of rem. size } x = \begin{cases} 0 & \text{if } r < x \\ x & \text{if } r \geq x \end{cases}$$

# From $r$ -work to number of jobs $N$

**Goal:** integral =  $N$



**Suffices:** integral = 1

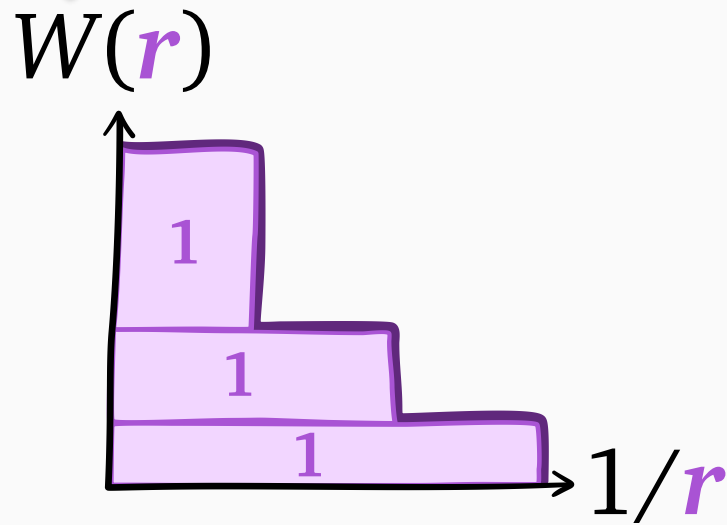


**Theorem:**

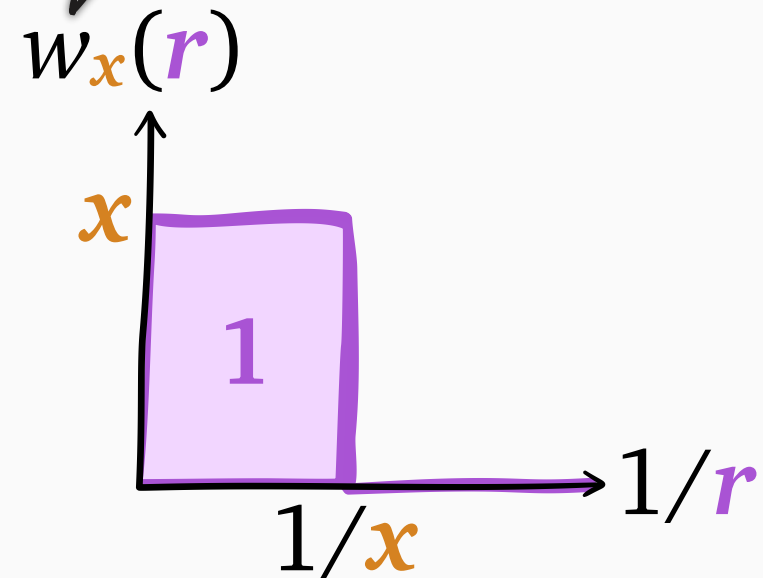
$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

# From $r$ -work to number of jobs $N$

Goal: integral =  $N$



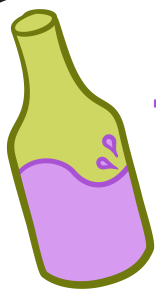
Suffices: integral = 1



**NEW!**

Theorem:

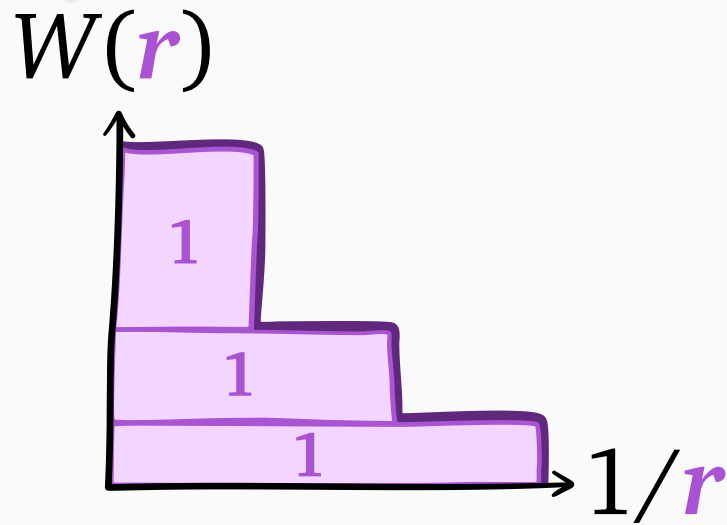
$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$



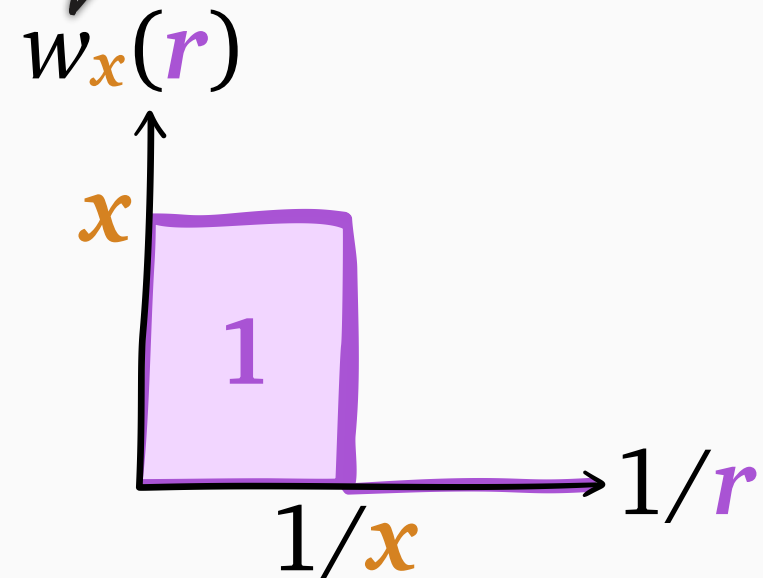
**WINE**

# From $r$ -work to number of jobs $N$

Goal: integral =  $N$



Suffices: integral = 1

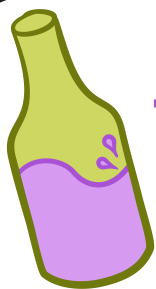


**NEW!**

Theorem:

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

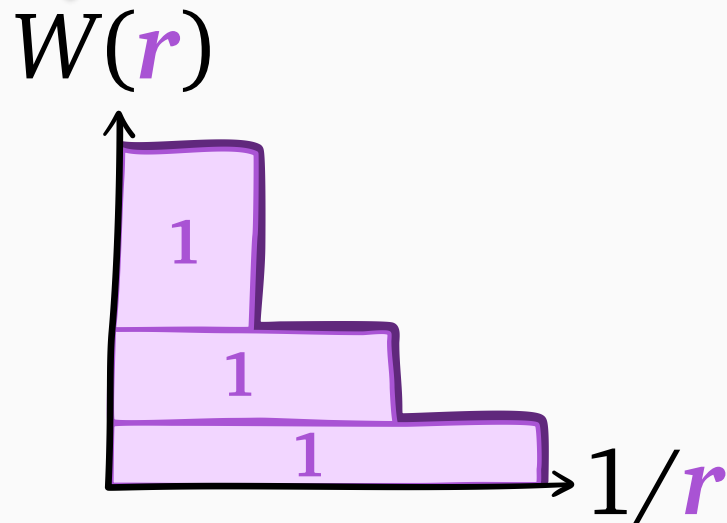
uses **rank** = rem. size



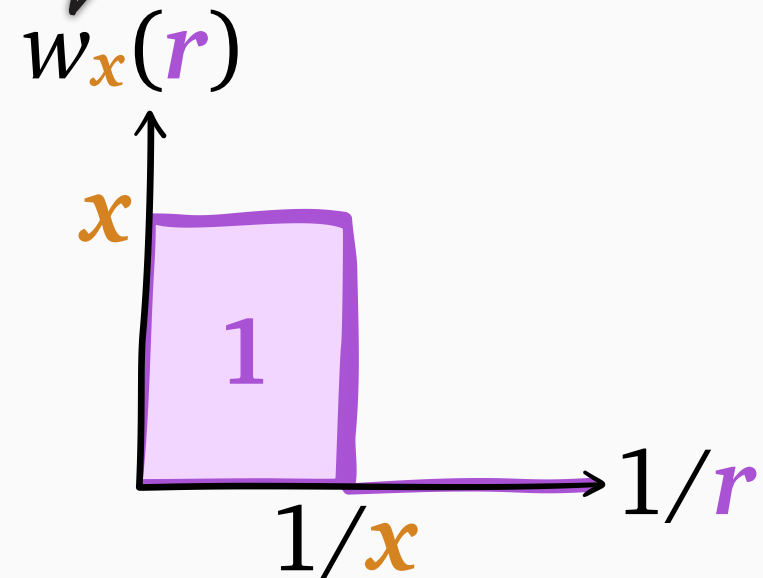
**WINE**

# From $r$ -work to number of jobs $N$

Goal: integral =  $N$



Suffices: integral = 1

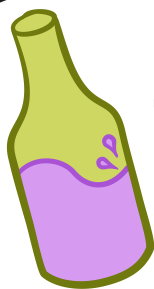


**NEW!**

Theorem: under *any* policy,

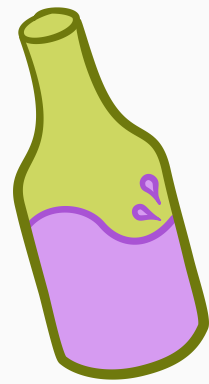
$$N = \int_0^{\infty} \frac{W(r)}{r^2} dr$$

uses **rank** = rem. size

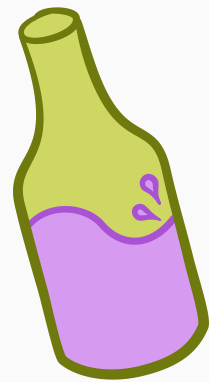


**WINE**





# Impact of WINE

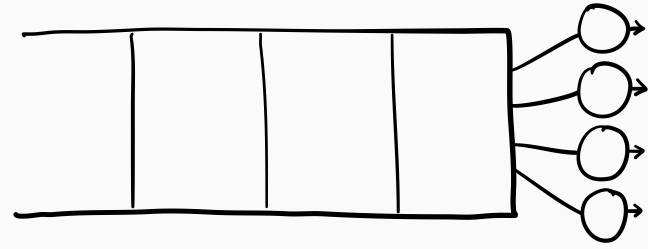


# Impact of WINE

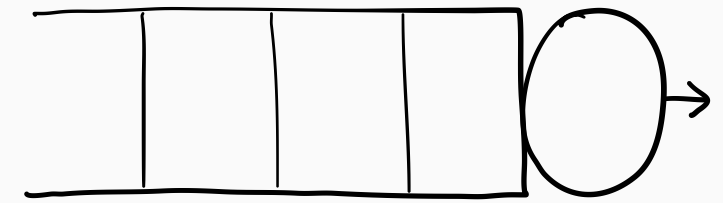


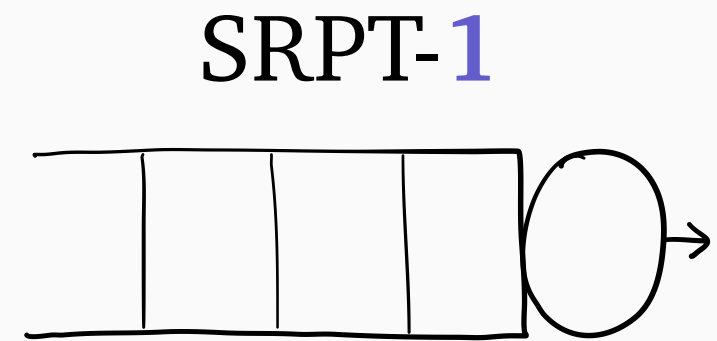
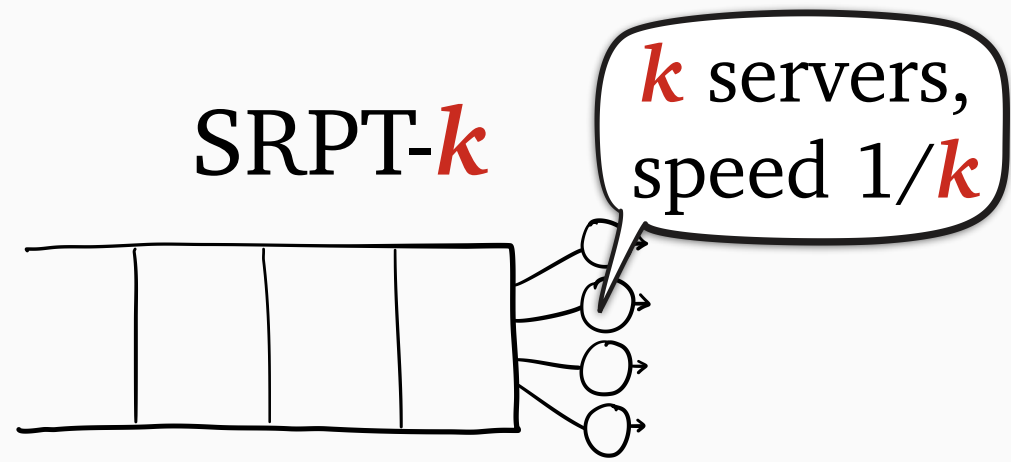
multiserver systems

SRPT-*k*

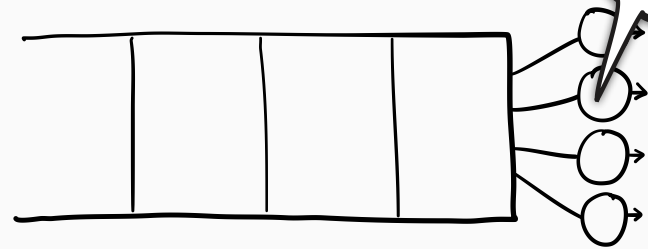


SRPT-1



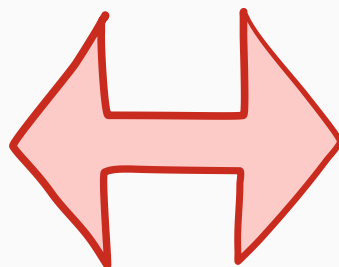
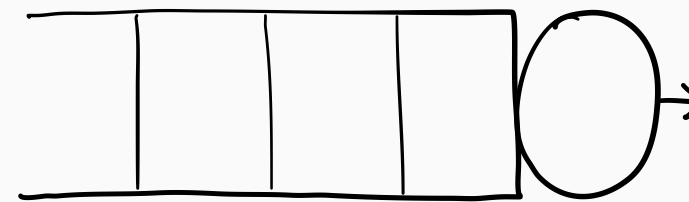


SRPT- $k$

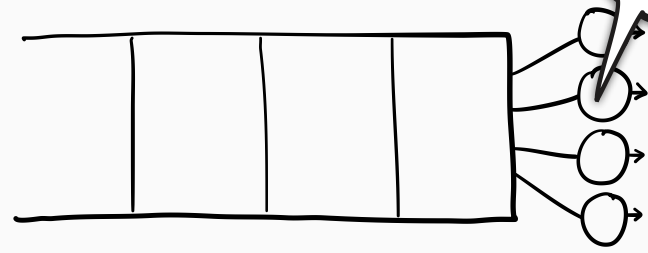


$k$  servers,  
speed  $1/k$

SRPT-1

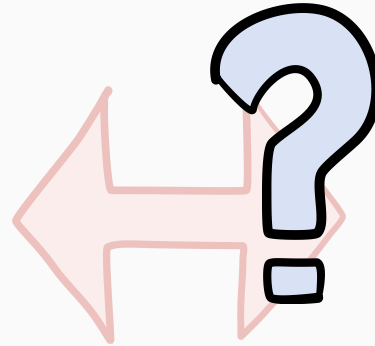
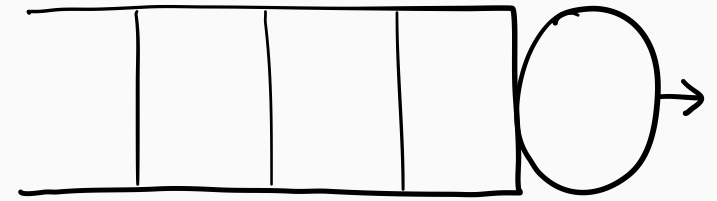


SRPT- $k$



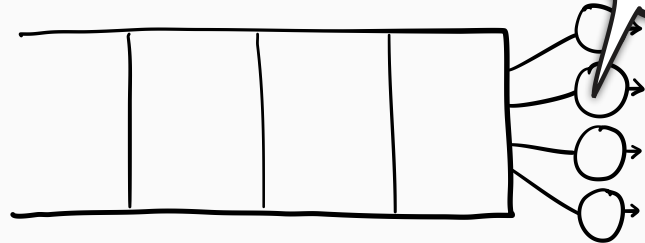
$k$  servers,  
speed  $1/k$

SRPT-1

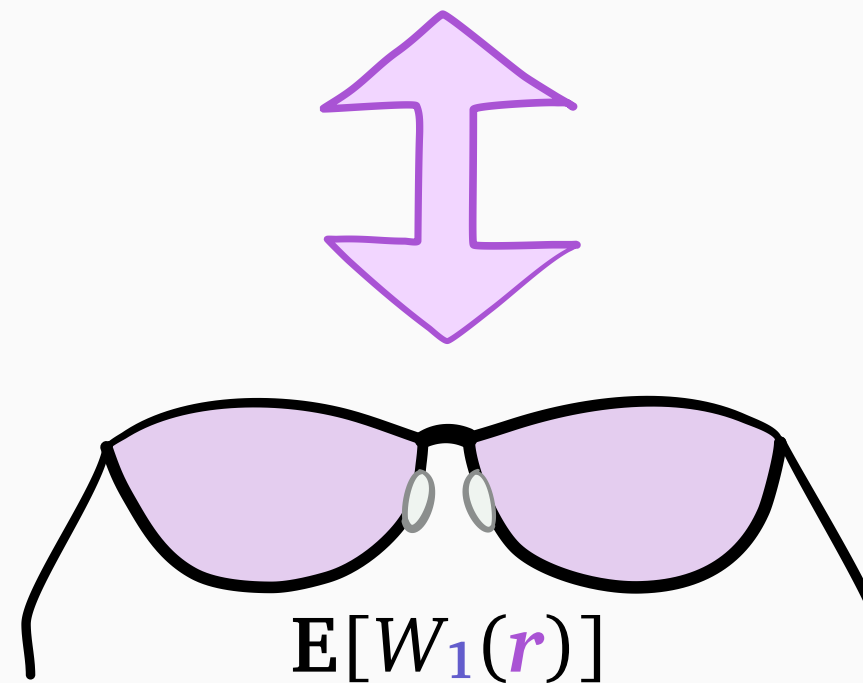
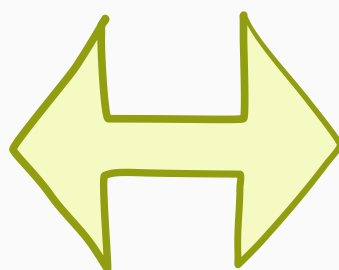
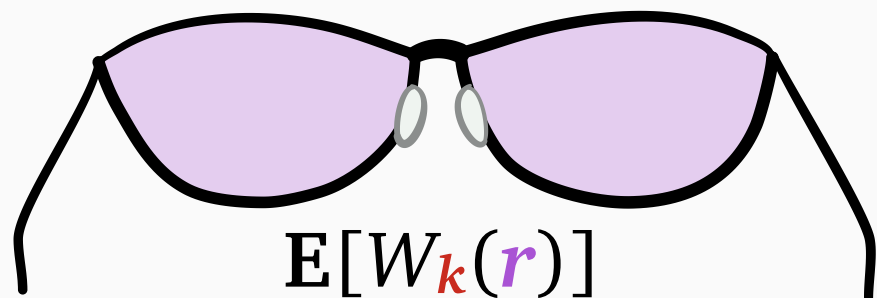
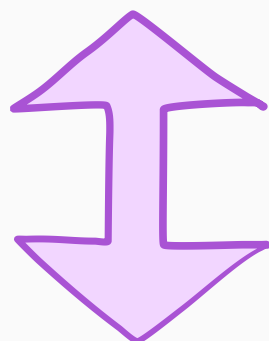
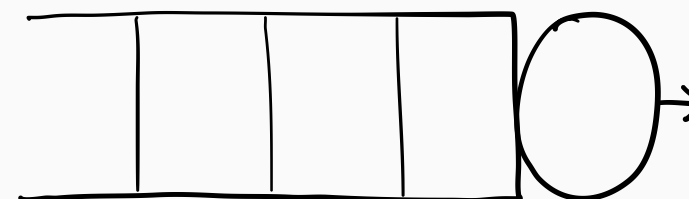


SRPT- $k$

$k$  servers,  
speed  $1/k$

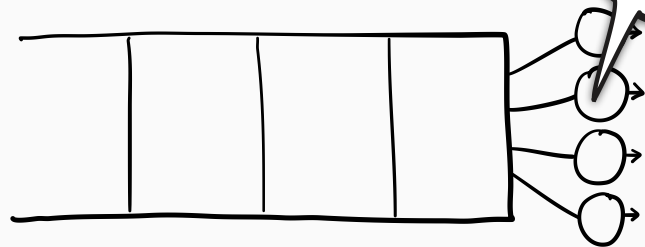


SRPT-1

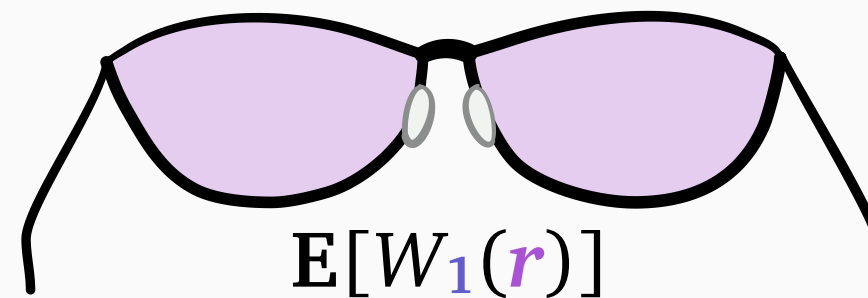
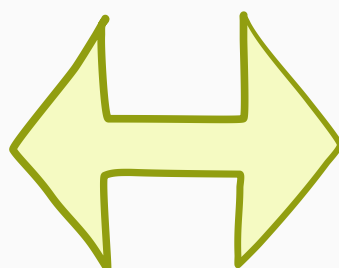
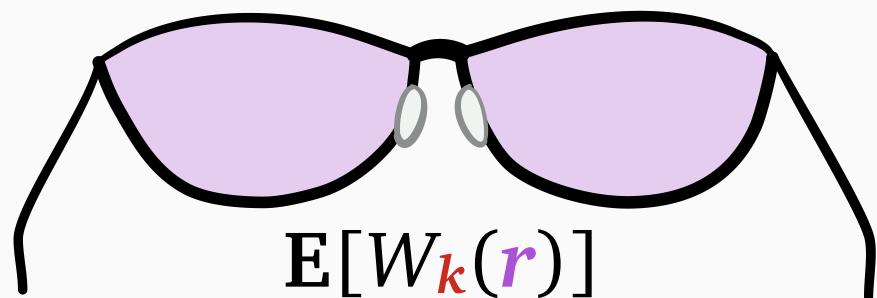
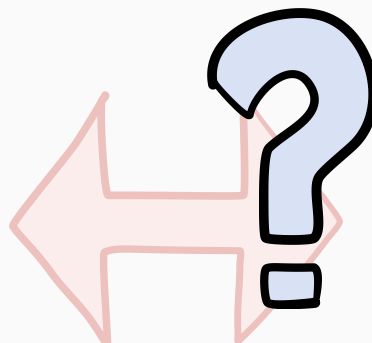
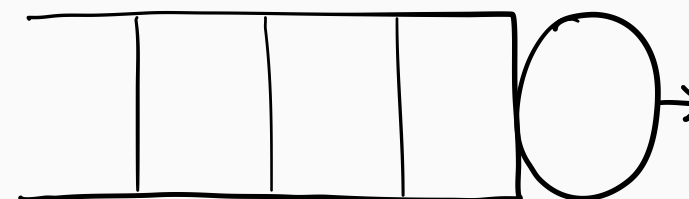


SRPT- $k$

$k$  servers,  
speed  $1/k$



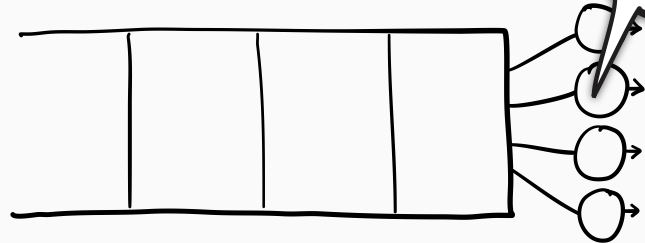
SRPT-1



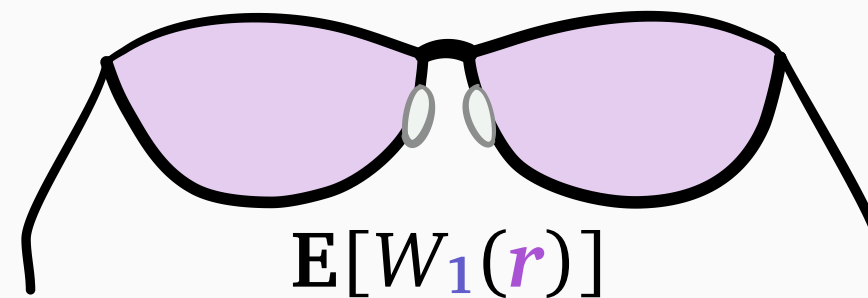
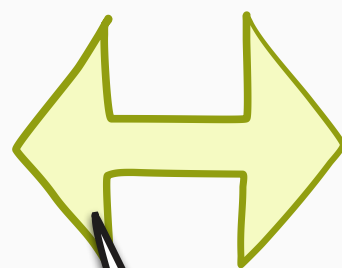
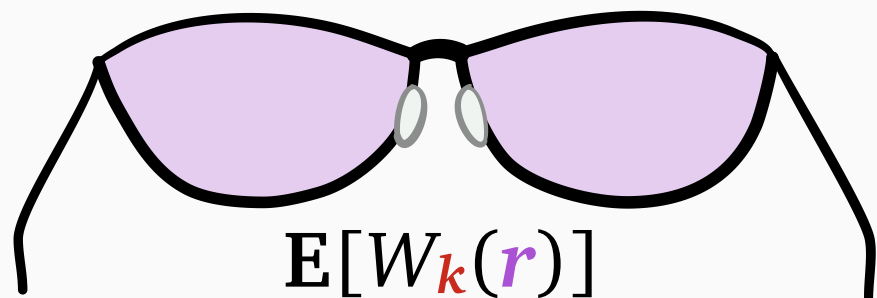
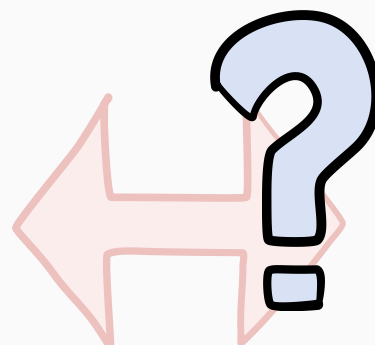
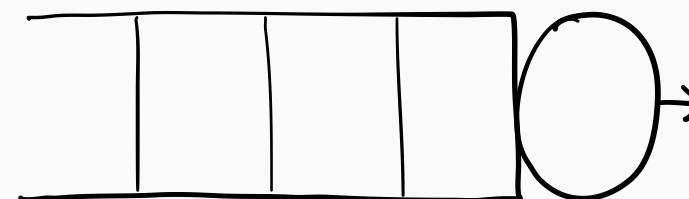


SRPT- $k$

$k$  servers,  
speed  $1/k$



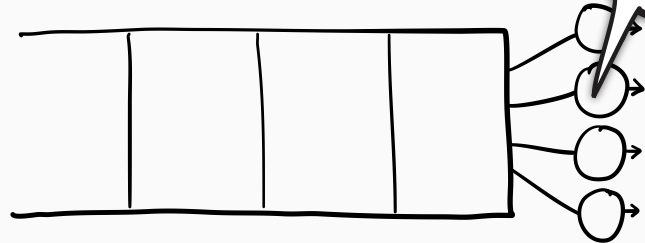
SRPT-1



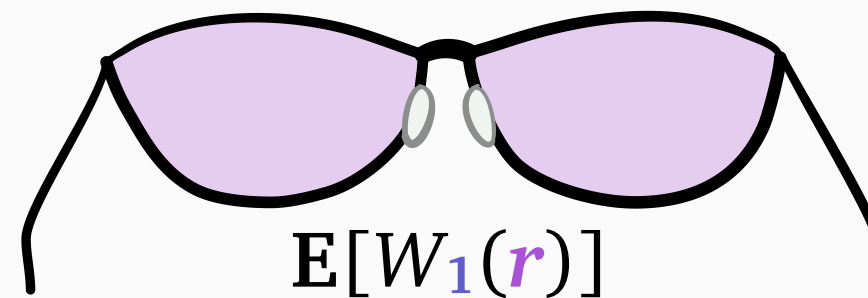
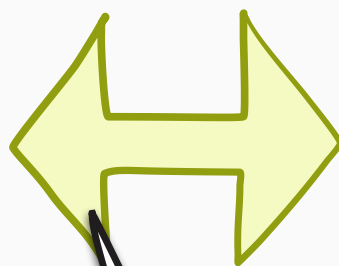
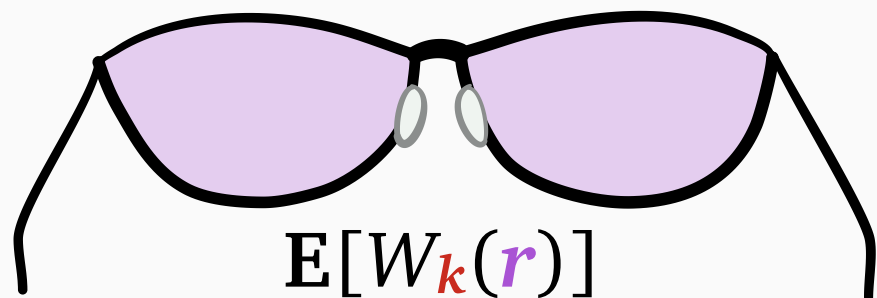
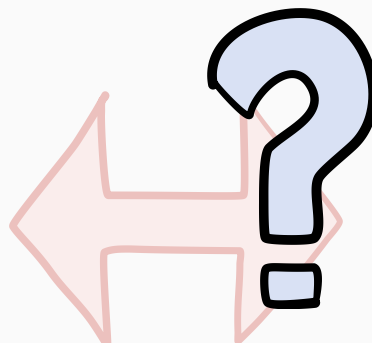
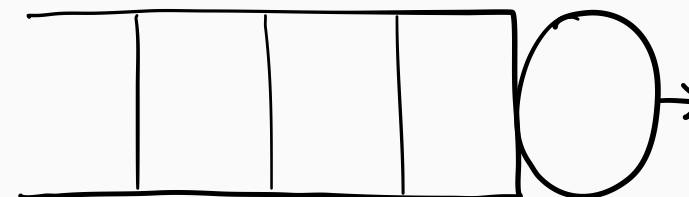
Lemma:  $r$ -work  
decomposition

SRPT- $k$

$k$  servers,  
speed  $1/k$



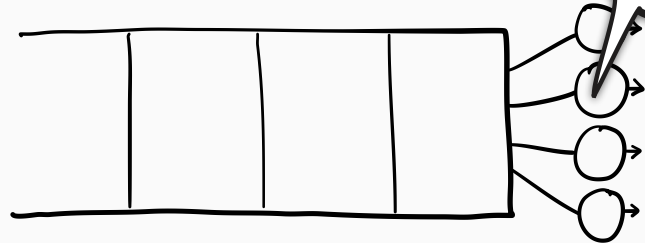
SRPT-1



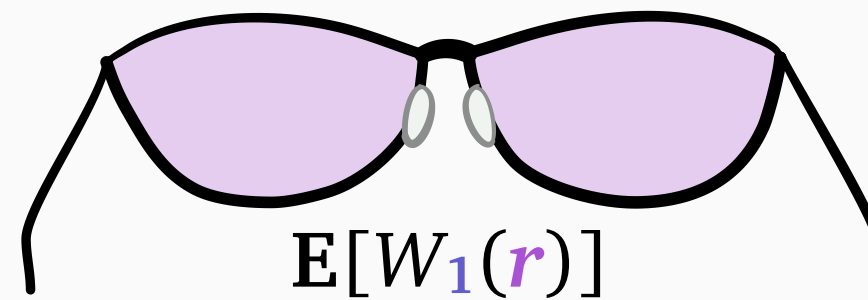
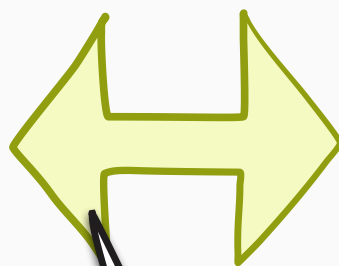
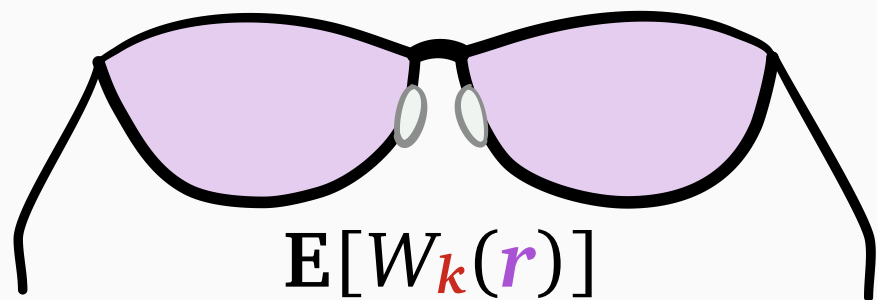
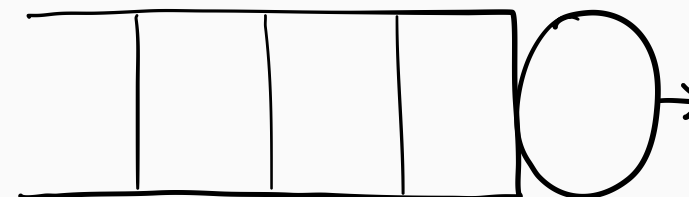
Lemma:  $r$ -work  
decomposition

SRPT- $k$

$k$  servers,  
speed  $1/k$



SRPT-1



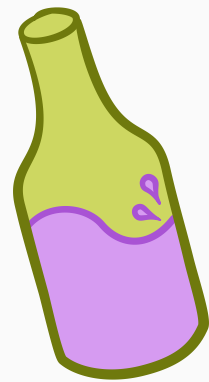
Lemma:  $r$ -work  
decomposition



# Impact of WINE



multiserver systems



# Impact of WINE

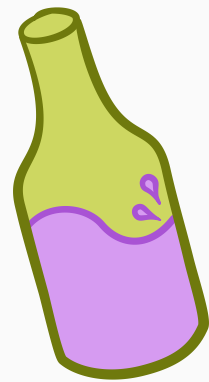


multiserver systems



*First analysis of*

- SRPT-*k*
- Gittins-*k*



# Impact of WINE



multiserver systems

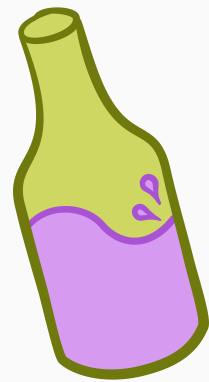


noisy size estimates



*First analysis of*

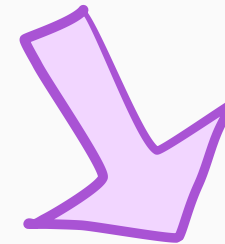
- SRPT-*k*
- Gittins-*k*



# Impact of WINE



multiserver systems



noisy size estimates



*First analysis of*

- SRPT-*k*
- Gittins-*k*

compare **noisy-info** *r*-work  
to **perfect-info** *r*-work

# WINE References

- [7] Grosf, Scully, and Harchol-Balter (2018). “SRPT for Multiserver Systems.” *Perform. Eval.* (PERFORMANCE 2018). **Winner: PERFORMANCE 2018 Best Student Paper Award.**
- [8] Grosf, Scully, and Harchol-Balter (2019). “Load Balancing Guardrails: Keeping Your Heavy Traffic on the Road to Low Response Times.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2019). **Winner: SIGMETRICS 2018 Outstanding Student Paper Award.**
- [9] Scully, Grosf, and Harchol-Balter (2021). “Optimal Multiserver Scheduling with Unknown Job Sizes in Heavy Traffic.” *Perform. Eval.* (PERFORMANCE 2020 issue).
- [10] Scully, Grosf, and Harchol-Balter (2020). “The Gittins Policy is Nearly Optimal in the M/G/k under Extremely General Conditions.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2021). **Winner: 2022 INFORMS George Nicholson Student Paper Competition.**
- [11] Scully (2021). “Bounding Mean Slowdown in Multiserver Systems.” *SIGMETRICS Perform. Eval. Rev.* (MAMA 2021).
- [12] Scully, Grosf, and Mitzenmacher (2022). “Uniform Bounds for Scheduling with Job Size Estimates.” *13th Innovations in Theoretical Computer Science Conference* (ITCS 2022).



# WINE References

[7] Grosf, Scully, and Harchol-Balter (2018). “SRPT for Multiserver Systems.” *Perform. Eval.* (PERFORMANCE 2018). **Winner: PERFORMANCE 2018 Best Student Paper Award.**

[8] Grosf, Scully, and Harchol-Balter (2019). “Load Balancing Guardrails: Keeping Your Heavy Traffic on the Road to Low Response Times.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2019). **Winner: SIGMETRICS 2018 Outstanding Student Paper Award.**

[9] Scully, Grosf, and Harchol-Balter (2021). “Optimal Multiserver Scheduling with Unknown Job Sizes in Heavy Traffic.” *Perform. Eval.* (PERFORMANCE 2020 issue).

[10] Scully, Grosf, and Harchol-Balter (2020). “The Gittins Policy is Nearly Optimal in the M/G/k under Extremely General Conditions.” *Proc. ACM Meas. Anal. Comput. Syst.* (SIGMETRICS 2021). **Winner: 2022 INFORMS George Nicholson Student Paper Competition.**

[11] Scully (2021). “Bounding Mean Slowdown in Multiserver Systems.” *SIGMETRICS Perform. Eval. Rev.* (MAMA 2021).

[12] Scully, Grosf, and Mitzenmacher (2022). “Uniform Bounds for Scheduling with Job Size Estimates.” *13th Innovations in Theoretical Computer Science Conference* (ITCS 2022).



# Outline: two new tools



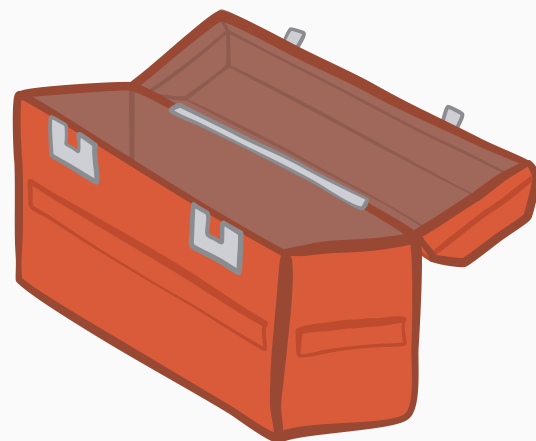
new unifying theory of  
single-server scheduling

- unknown sizes
- preemption limitations



new queueing identity to  
complement Little's Law

- multiserver systems
- noisy size estimates



# Outline: two new tools



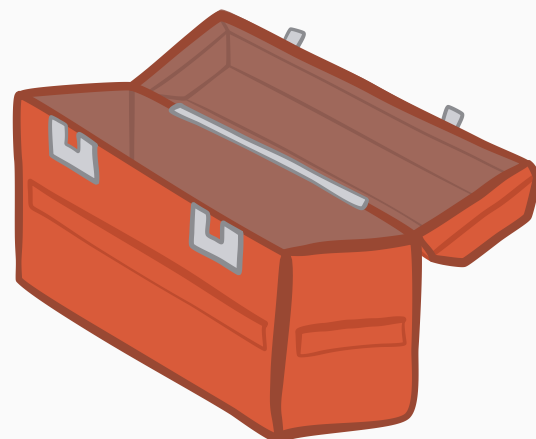
new unifying theory of single-server scheduling

- unknown sizes
- preemption limitations



new queueing identity to complement Little's Law

- multiserver systems
- noisy size estimates







*Thesis statement?*



*Thesis statement?*

Yes, queueing theory can!

**Past:**  
Analyze policies one-by-one

**Past:**

Analyze policies one-by-one



**Present:**

Analyze all **rank**-function-based policies at once



MLPS, SMART

**Past:**

Analyze policies one-by-one



**Present:**

Analyze all **rank**-function-based policies at once

MLPS, SMART

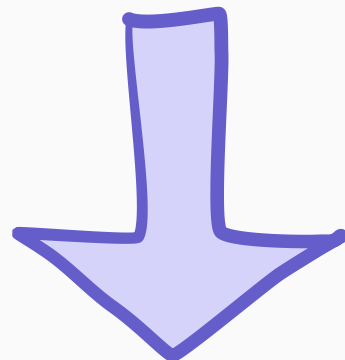
**Past:**

Analyze policies one-by-one



**Present:**

Analyze all **rank**-function-based policies at once



**Future:**

Other policy classes?

**Past:**

Scheduling in  $M/G/1$ , only FCFS in  $M/G/k$

**Past:**

Scheduling in  $M/G/1$ , only FCFS in  $M/G/k$



**Present:**

Scheduling in  $M/G/k$

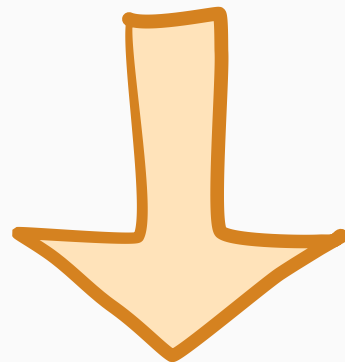
**Past:**

Scheduling in  $M/G/1$ , only FCFS in  $M/G/k$



**Present:**

Scheduling in  $M/G/k$

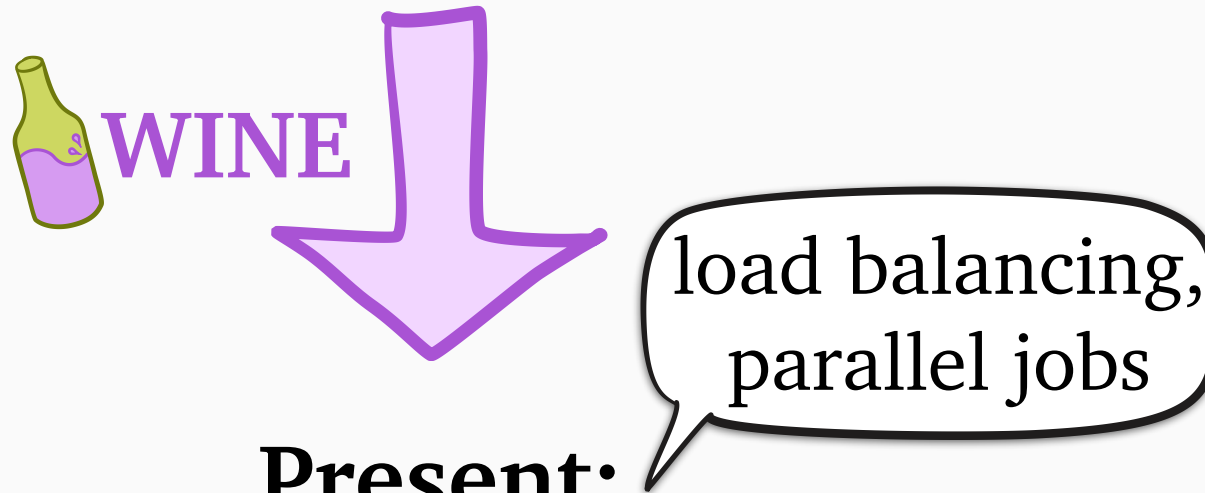


**Future:**

Scheduling in complex multiserver architectures?

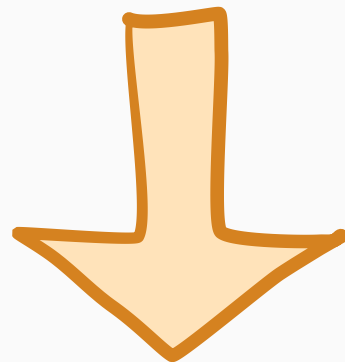
## Past:

Scheduling in  $M/G/1$ , only FCFS in  $M/G/k$



## Present:

Scheduling in  $M/G/k$



## Future:

Scheduling in complex multiserver architectures?

**Past:**

**Gittins** for specific uncertainty models, one-by-one

**Past:**

**Gittins** for specific uncertainty models, one-by-one



**Present:**

**Gittins** for any service-based uncertainty model



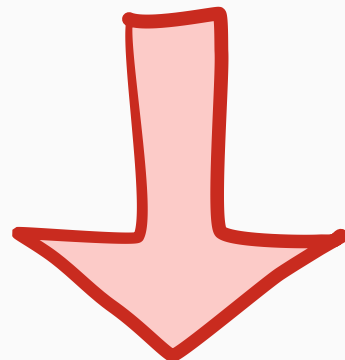
**Past:**

**Gittins** for specific uncertainty models, one-by-one



**Present:**

**Gittins** for any service-based uncertainty model



**Future:**

Underspecified models? “Restless” models?

# Yes, queueing theory can!



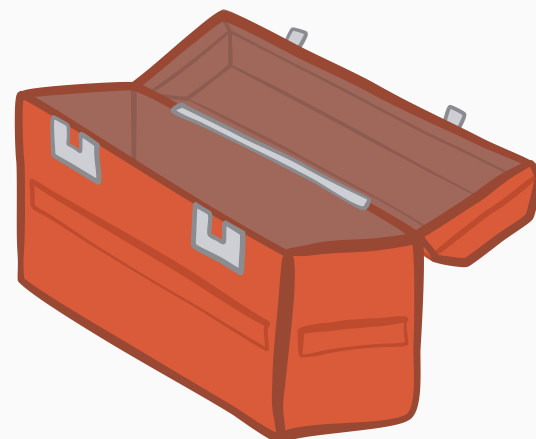
new unifying theory of  
single-server scheduling

- unknown sizes
- preemption limitations



new queueing identity to  
complement Little's Law

- multiserver systems
- noisy size estimates



**SOAP** def. and non-examples

Preemption checkpoints

Limited priority levels

Tail decay: heavy vs. light

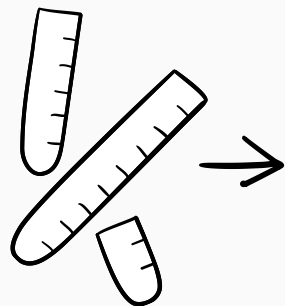
*r*-work decomposition

SRPT-*k* and Gittins-*k*  $E[T]$

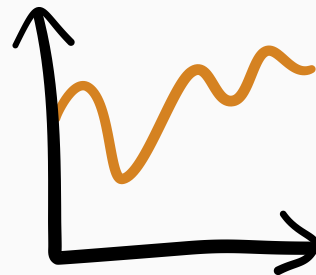
# SOAP

Schedule **O**rdered by **A**ge-based **P**riority

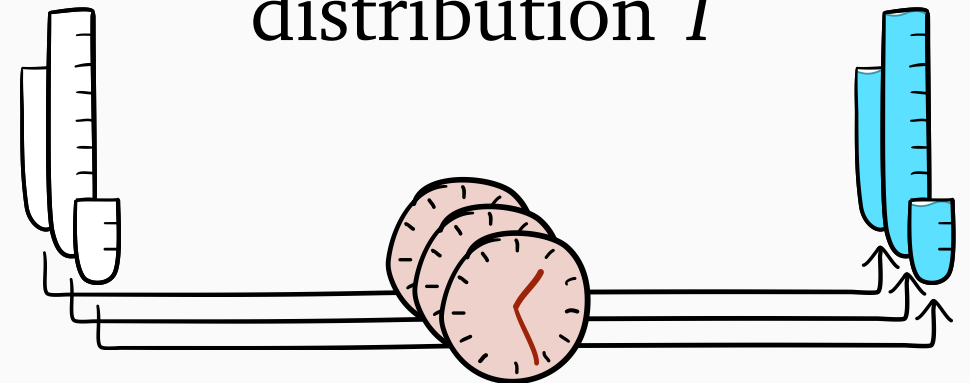
stochastic arrival  
process  $\lambda, S$



**rank** function



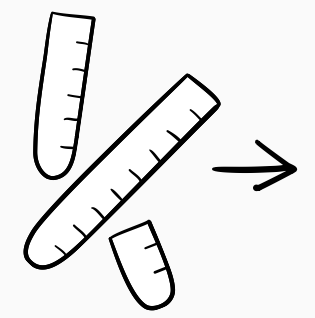
response time  
distribution  $T$



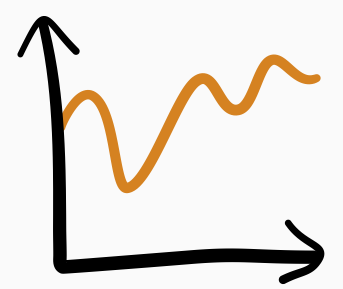
# SOAP

Schedule **O**rdered by **A**ge-based **P**riority

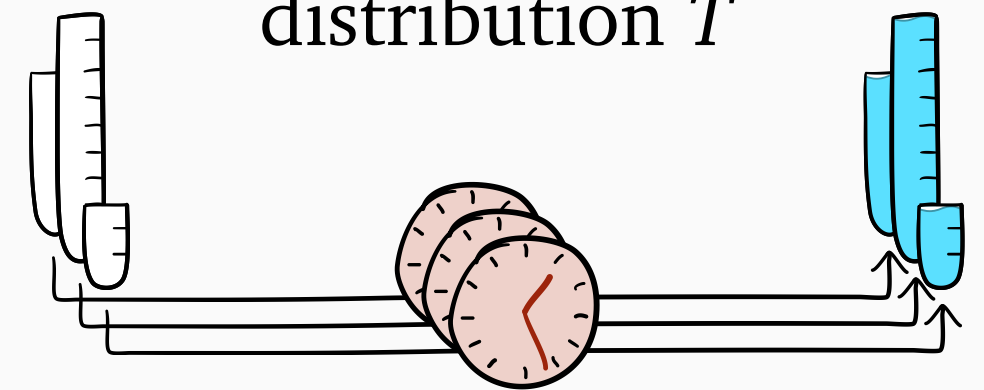
stochastic arrival process  $\lambda, S$



**rank** function



response time distribution  $T$



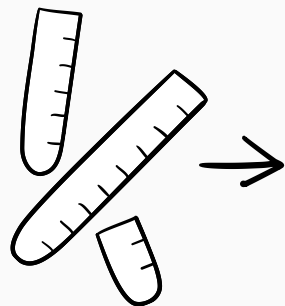
$(\text{label}, \text{age}) \mapsto \text{rank}$   
where labels are static

# SOAP

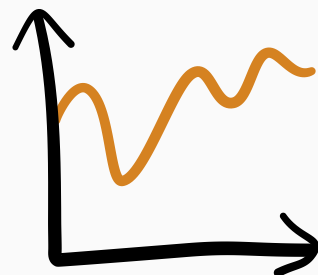
Schedule **O**rdered by **A**ge-based **P**riority

label dist.  $L$ ,  
 $(L, S)$  pairs i.i.d.

stochastic arrival  
process  $\lambda, S$

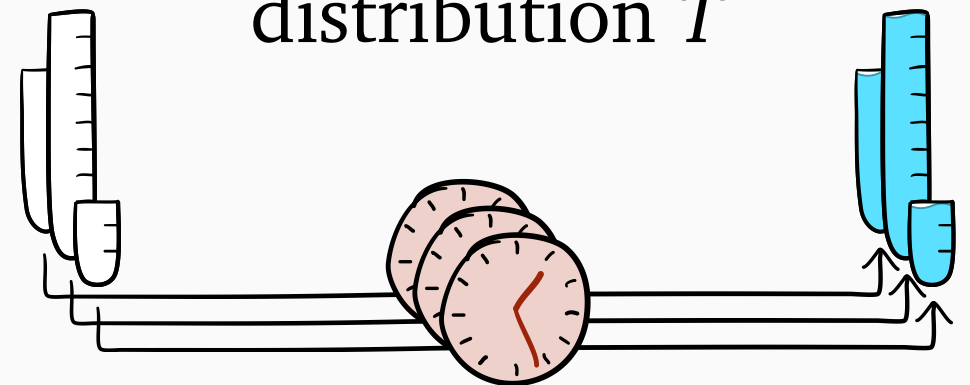


rank function



$(\text{label}, \text{age}) \mapsto \text{rank}$   
where labels are static

response time  
distribution  $T$



back

# What's *not* a **SOAP** policy?

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**



# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

Job's **rank** changes while it's not in service

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

Job's **rank** changes while it's not in service

*EDF: can't have **rank** = time until deadline*

earliest deadline first

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

Job's **rank** changes while it's not in service

*EDF: can't have **rank** = time until deadline*

Job's **rank** depends on a non-i.i.d. label

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

Job's **rank** changes while it's not in service

*EDF: can't have **rank** = time until deadline*

Job's **rank** depends on a non-i.i.d. label

*EDF: can't have **rank** = deadline*

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

Job's **rank** changes while it's not in service

*EDF: can't have **rank** = time until deadline*

Job's **rank** depends on a non-i.i.d. label

*EDF: can't have **rank** = deadline*

Tiebreaking not FCFS (or LCFS)

# What's *not* a **SOAP** policy?

Scheduling not based on each job having a **rank**

*Nudge: can't use arrival sequence adjacency*

Job's **rank** changes while it's not in service

*EDF: can't have **rank** = time until deadline*

Job's **rank** depends on a non-i.i.d. label

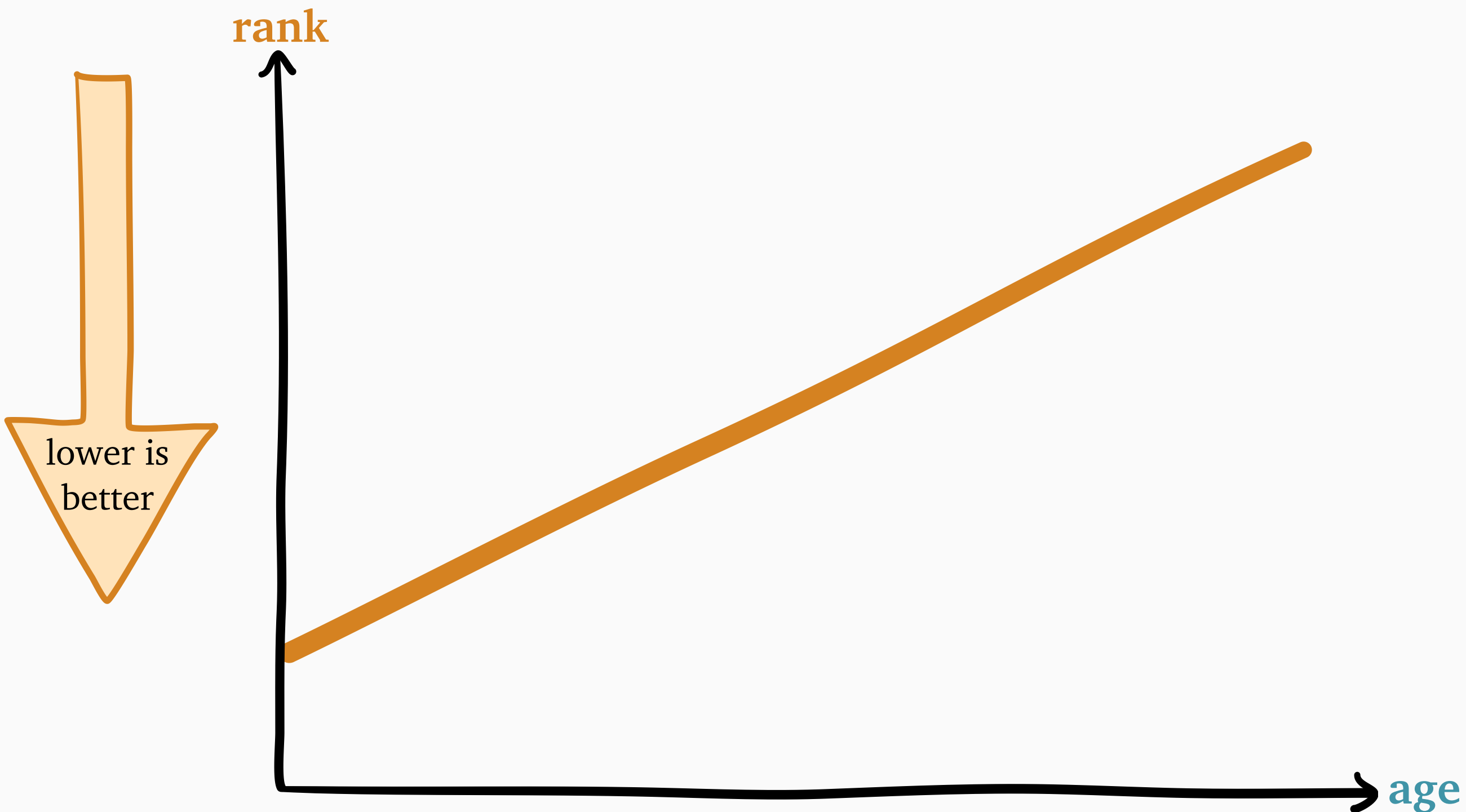
*EDF: can't have **rank** = deadline*

Tiebreaking not FCFS (or LCFS)

*PS: can't break ties by sharing (or randomizing)*

# Preemption checkpoints

Can only preempt only at *checkpoint ages*

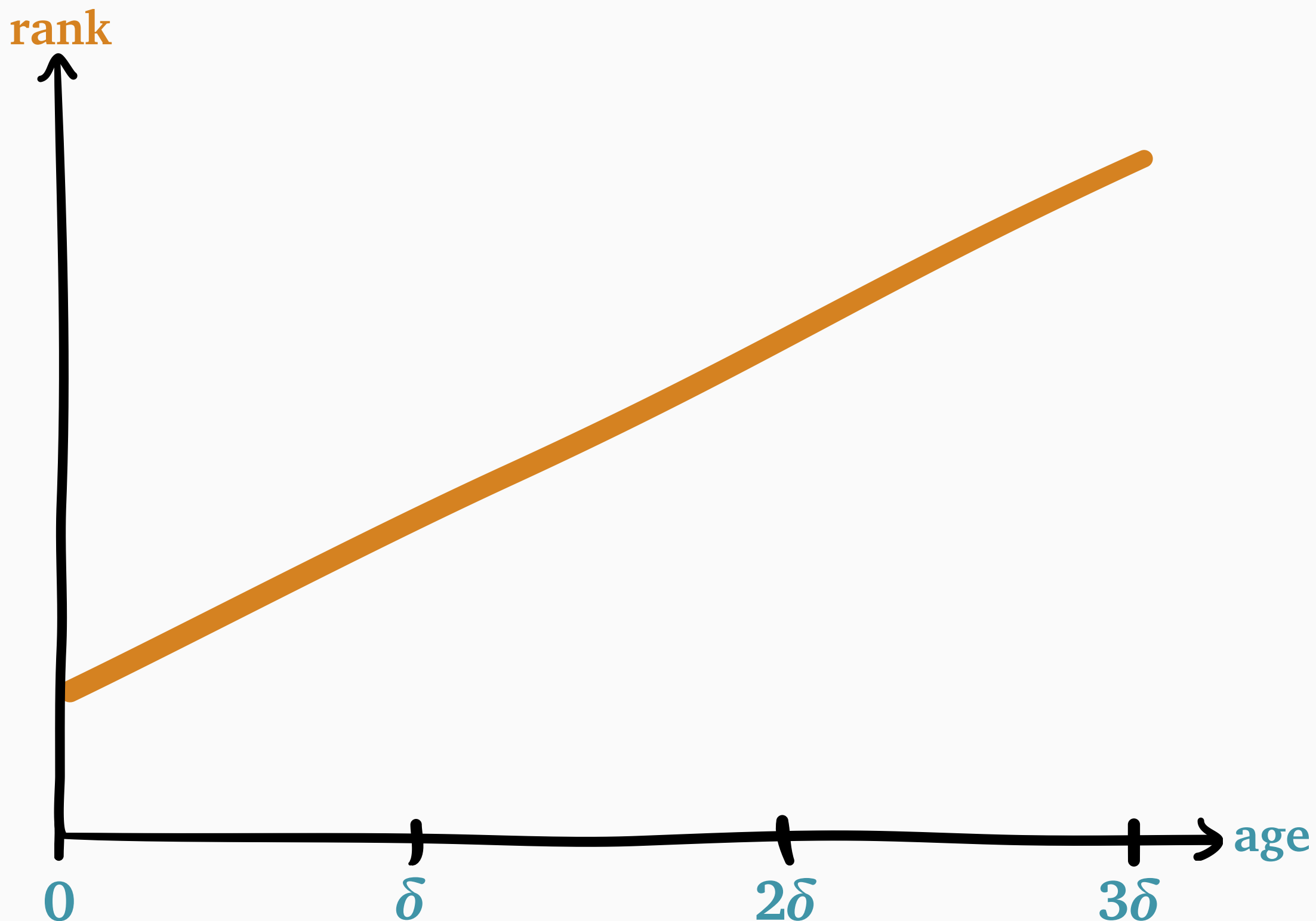




back

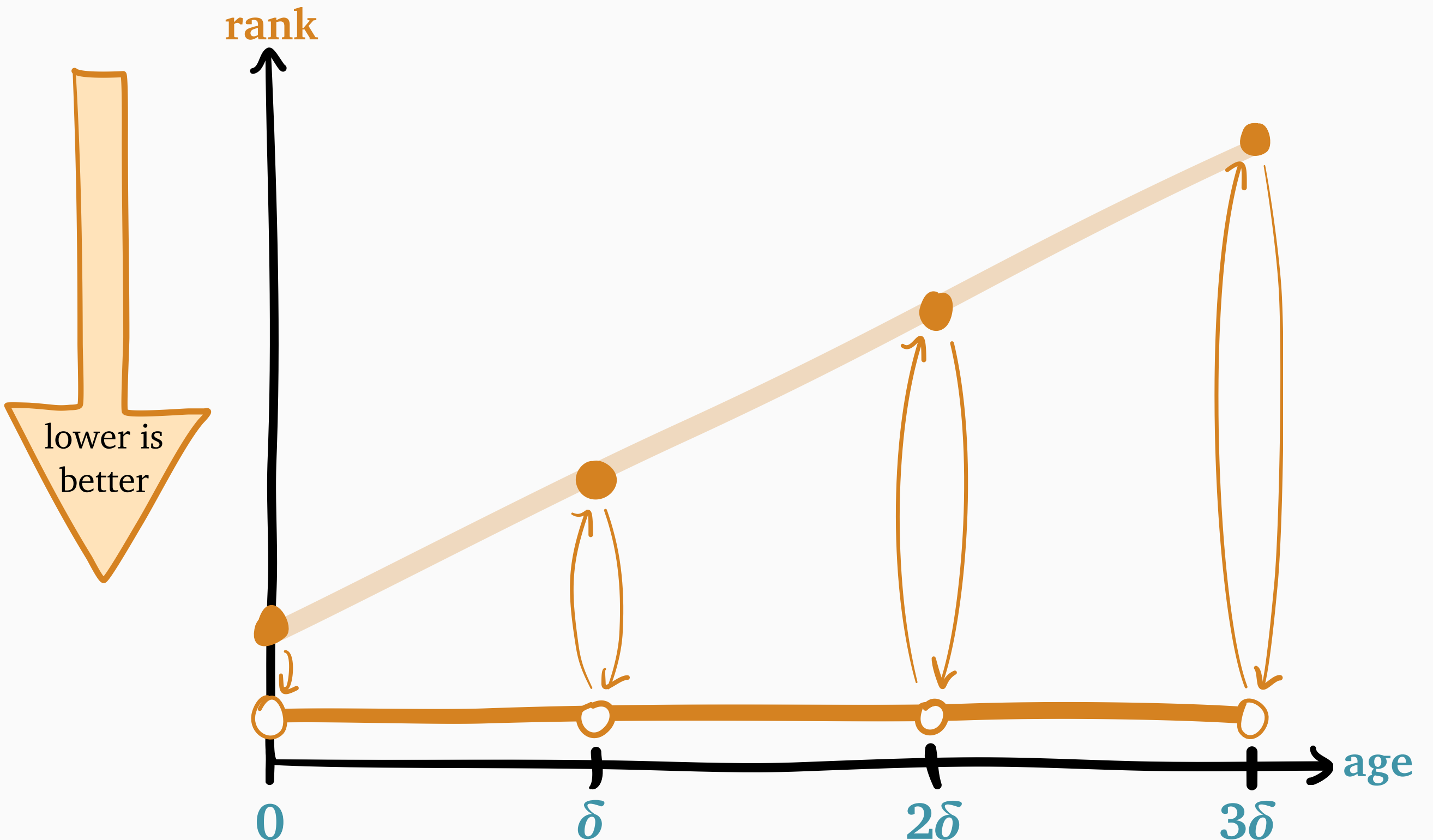
# Preemption checkpoints

Can only preempt only at *checkpoint ages*



# Preemption checkpoints

Can only preempt only at *checkpoint ages*



# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

optimal packet size?

# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

- large  $\delta$ : less overhead

optimal packet size?

# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

- large  $\delta$ : less overhead
- small  $\delta$ : more precise scheduling

optimal packet size?

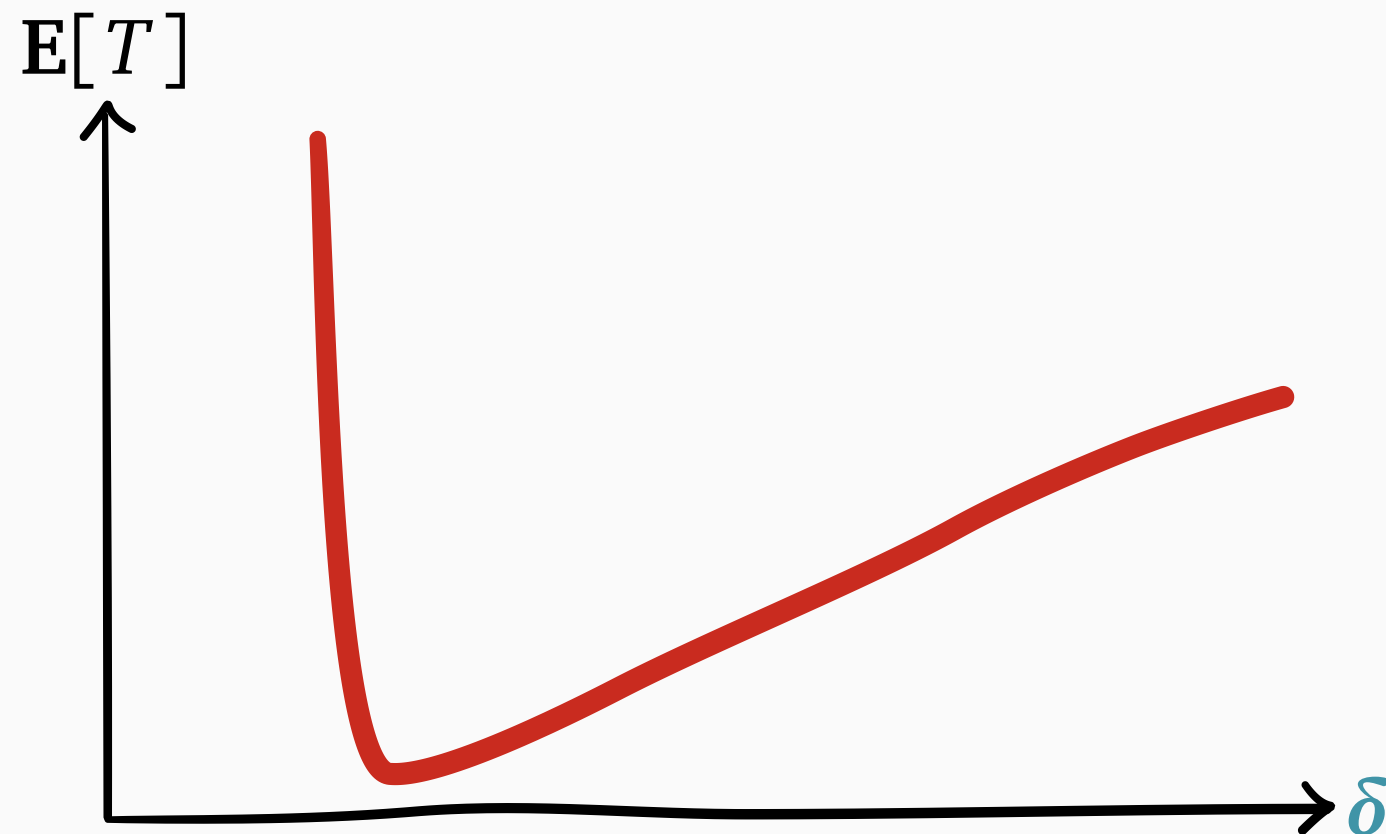
# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

- large  $\delta$ : less overhead
- small  $\delta$ : more precise scheduling

optimal packet size?





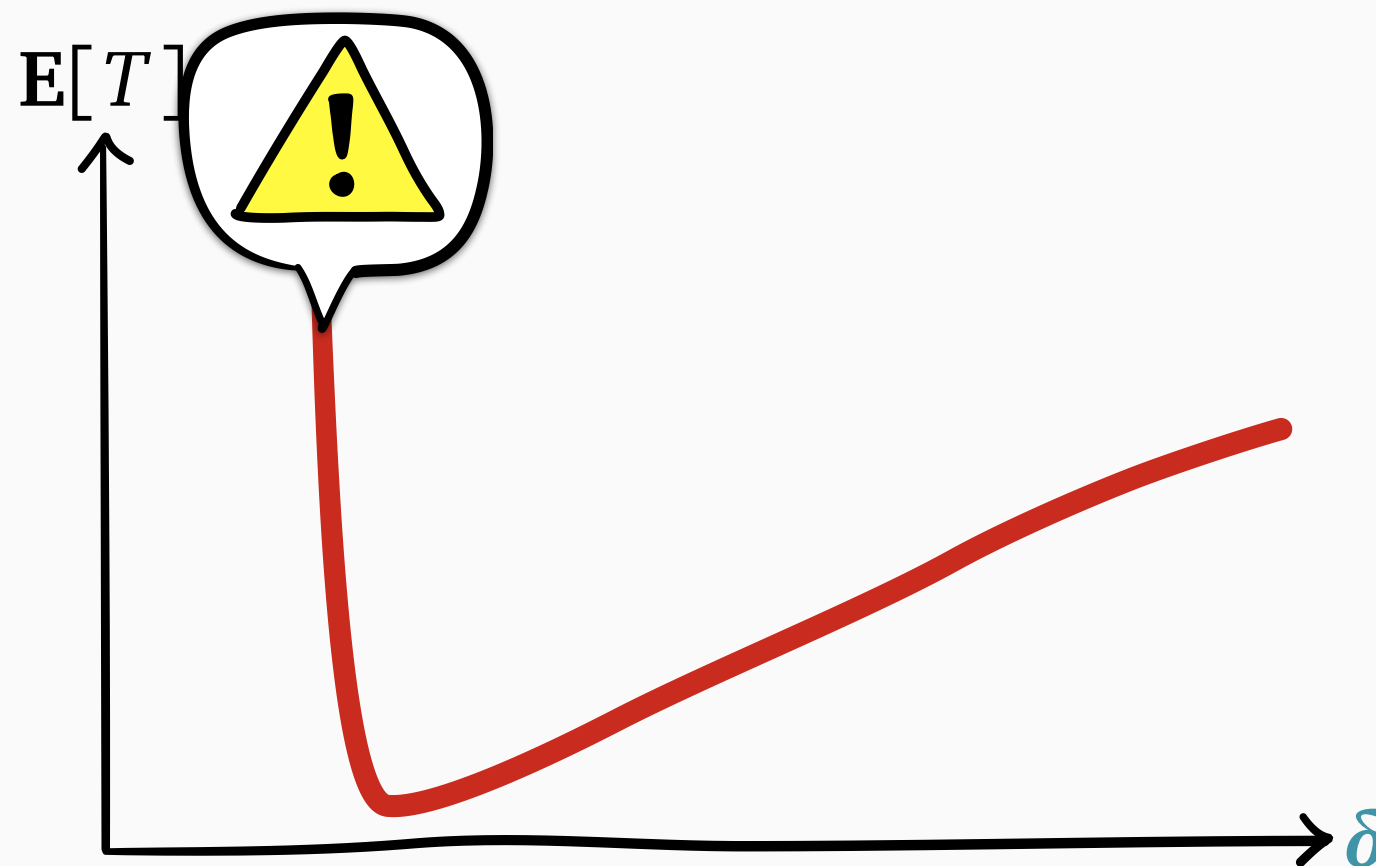
# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

- large  $\delta$ : less overhead
- small  $\delta$ : more precise scheduling

optimal packet size?



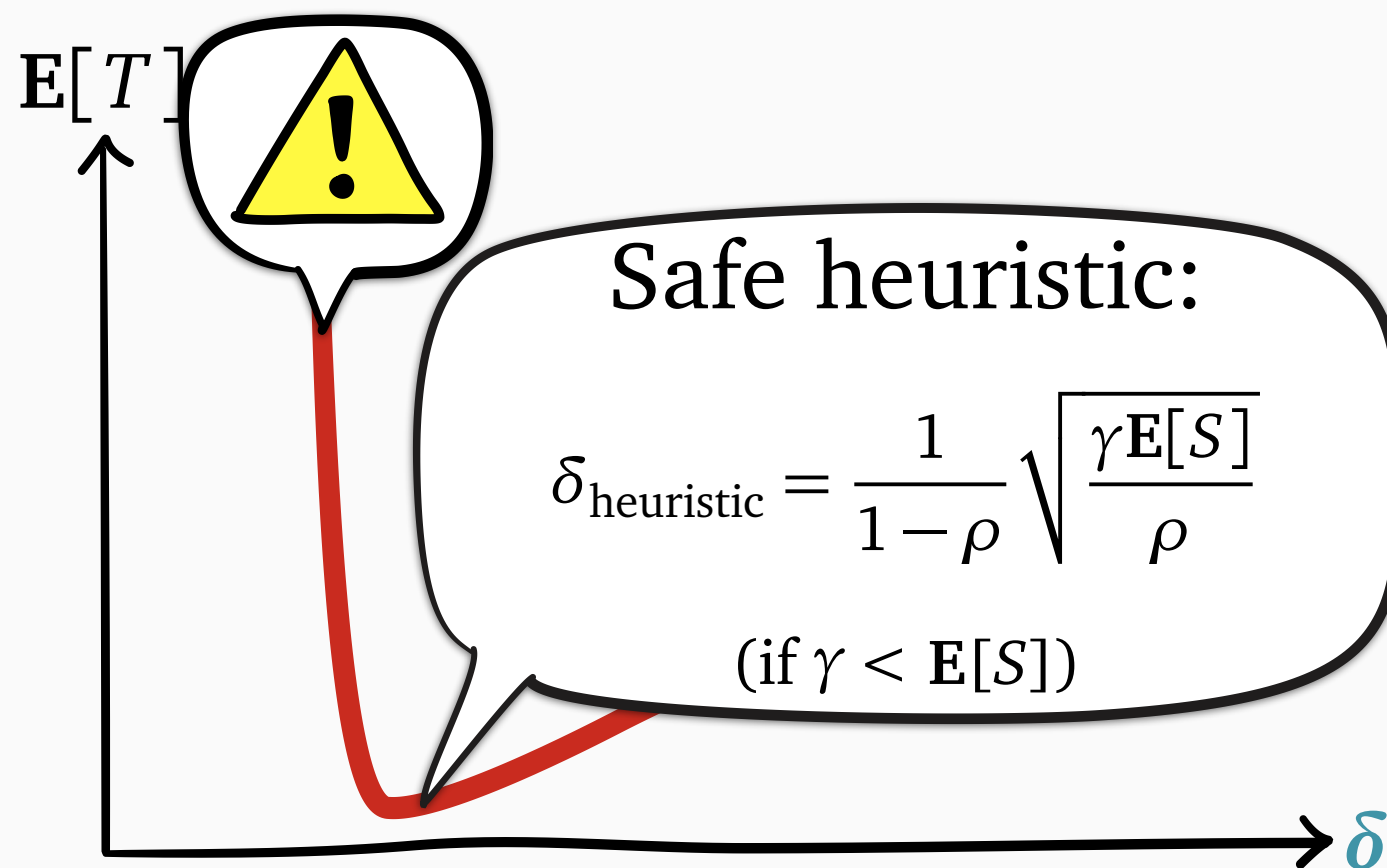
# Checkpoint frequency?

Suppose each checkpoint incurs an overhead  $\gamma$

What is the optimal gap  $\delta$  between checkpoints?

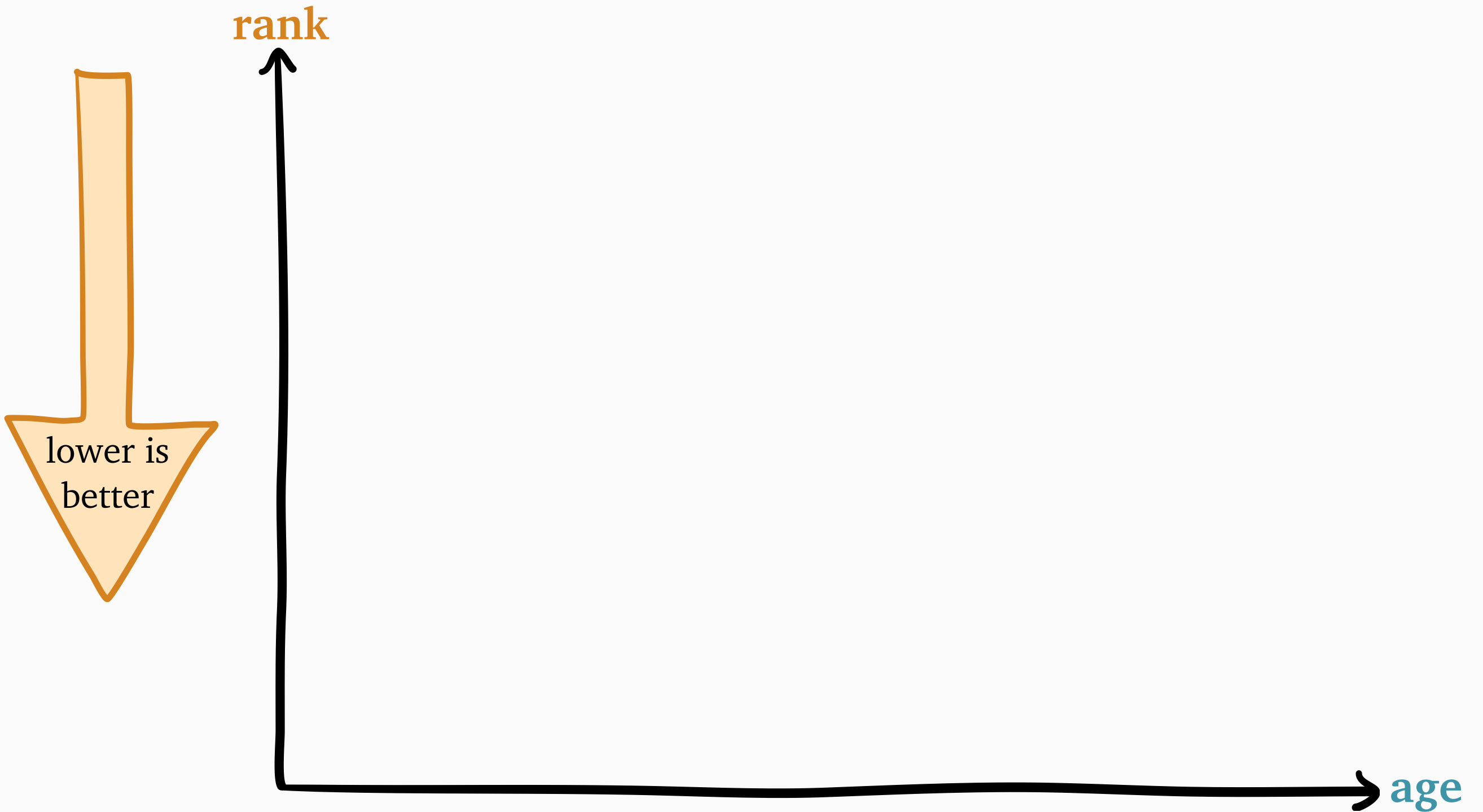
- large  $\delta$ : less overhead
- small  $\delta$ : more precise scheduling

optimal packet size?

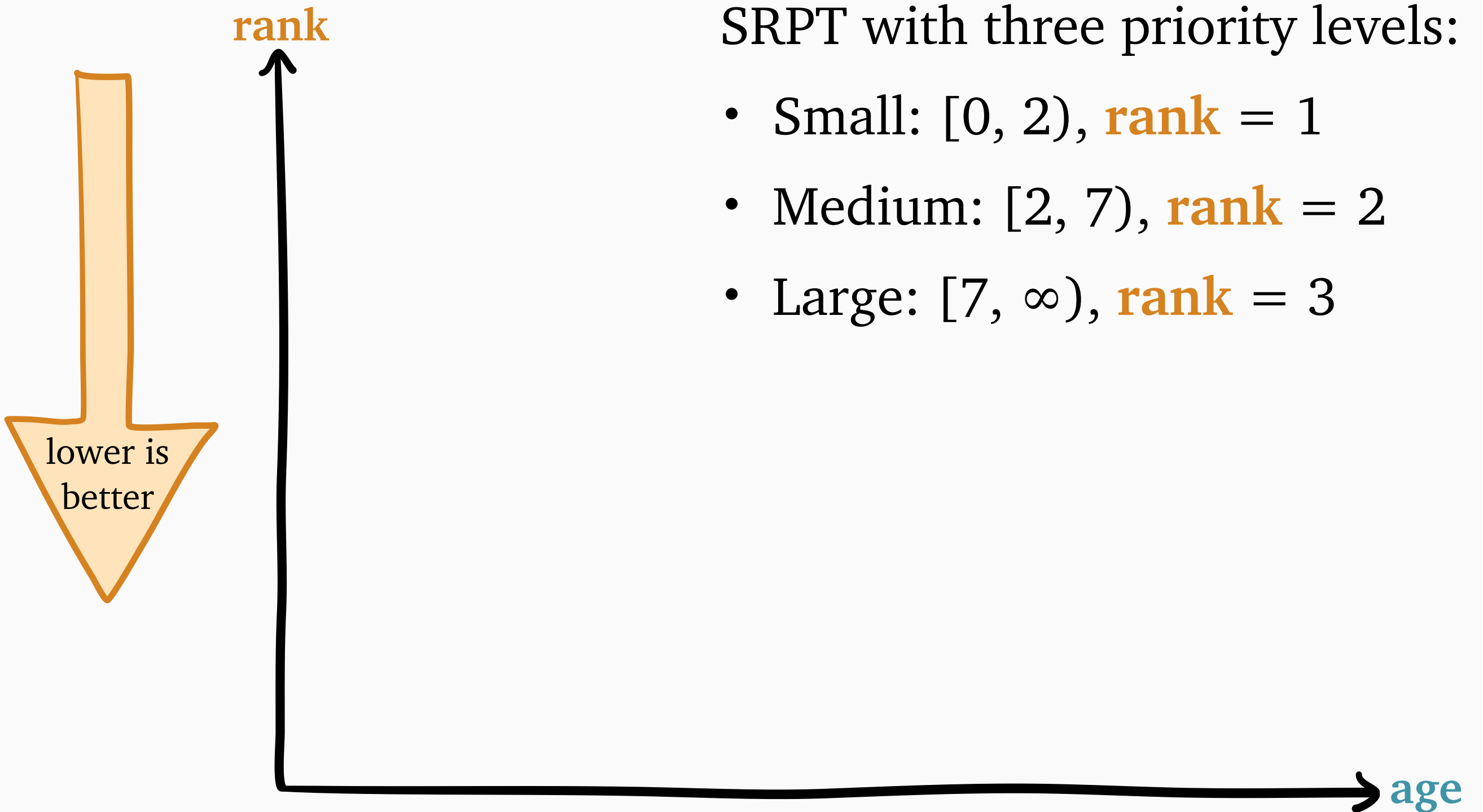


back

# Limited priority levels



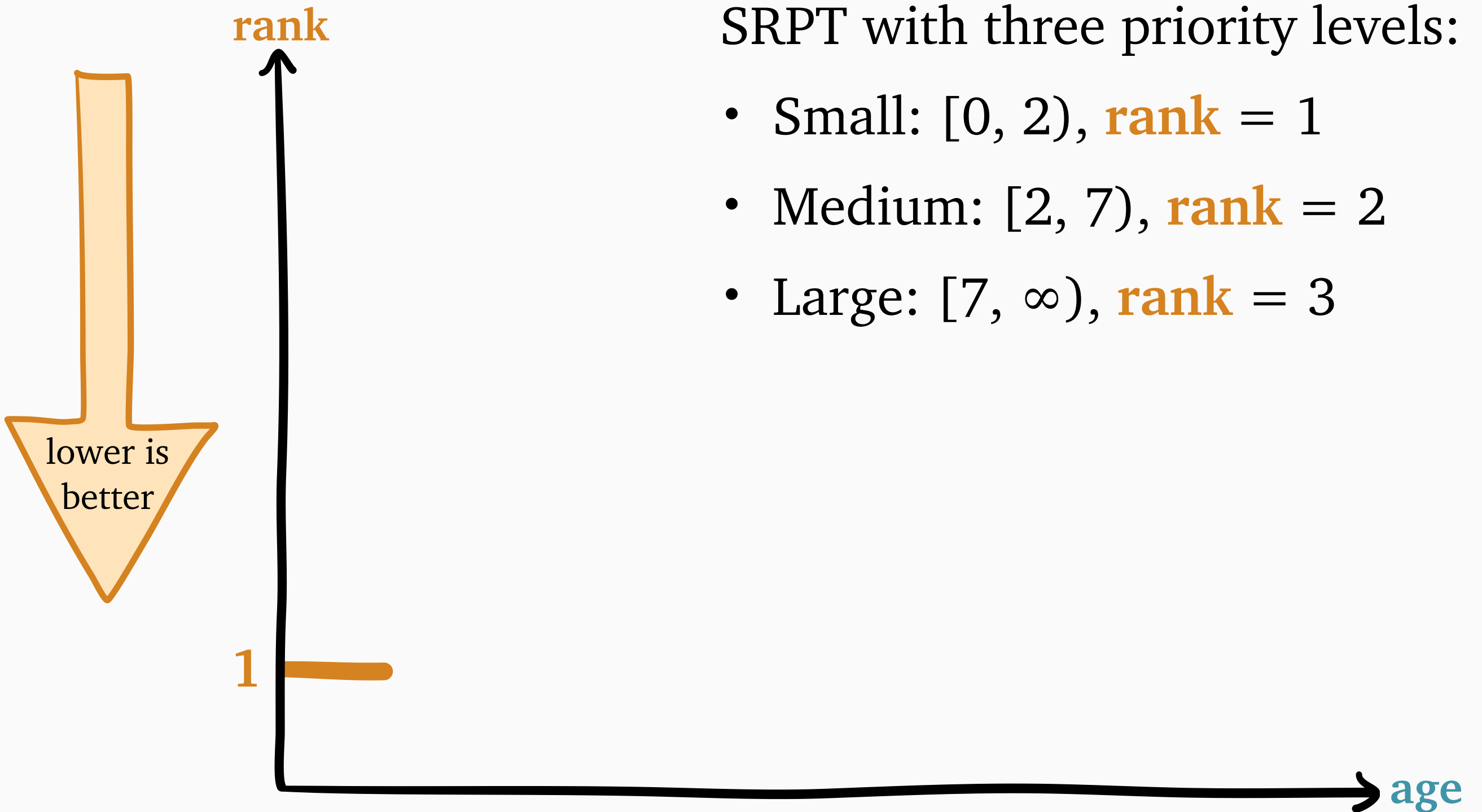
# Limited priority levels



# Limited priority levels

SRPT with three priority levels:

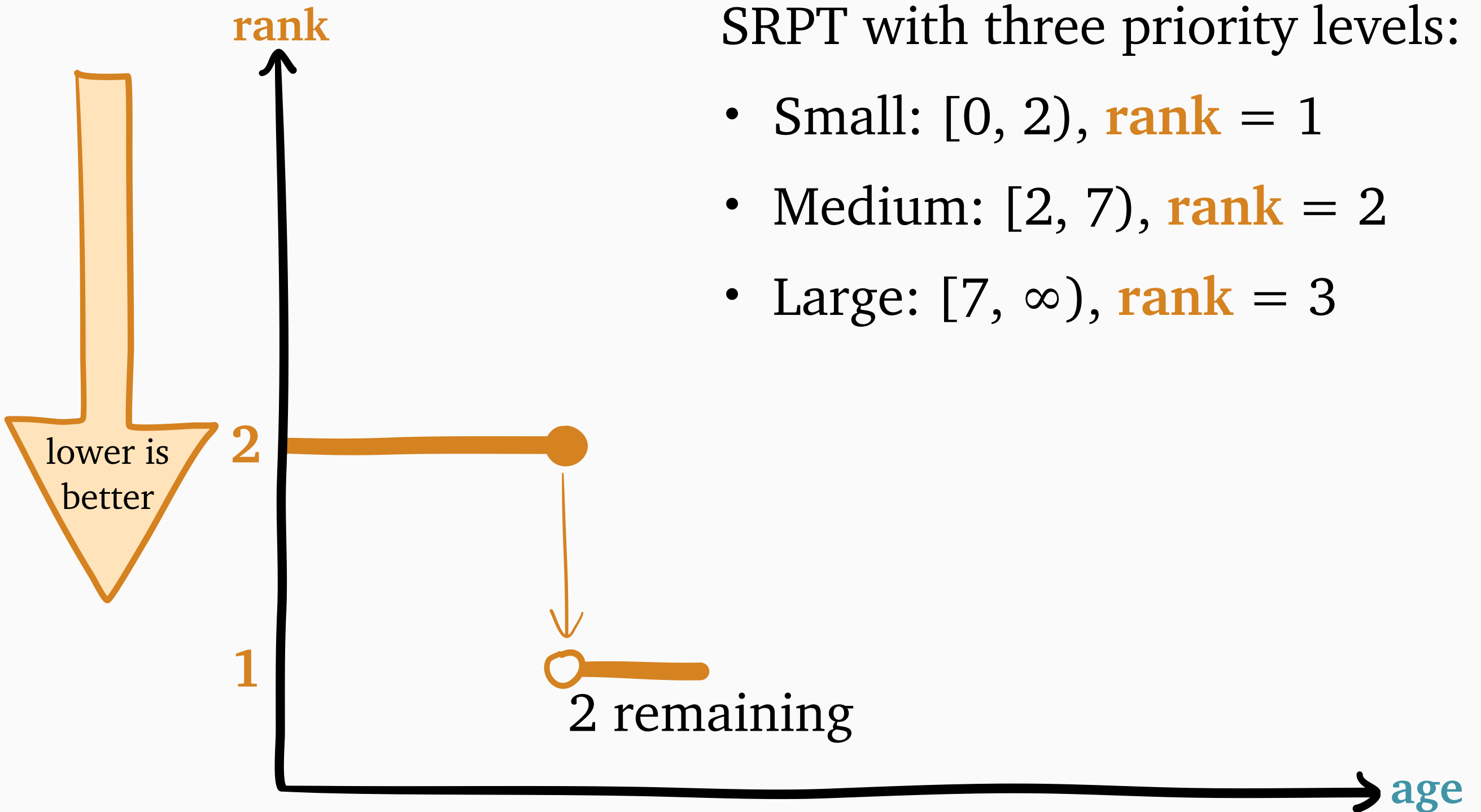
- Small:  $[0, 2)$ , **rank** = 1
- Medium:  $[2, 7)$ , **rank** = 2
- Large:  $[7, \infty)$ , **rank** = 3



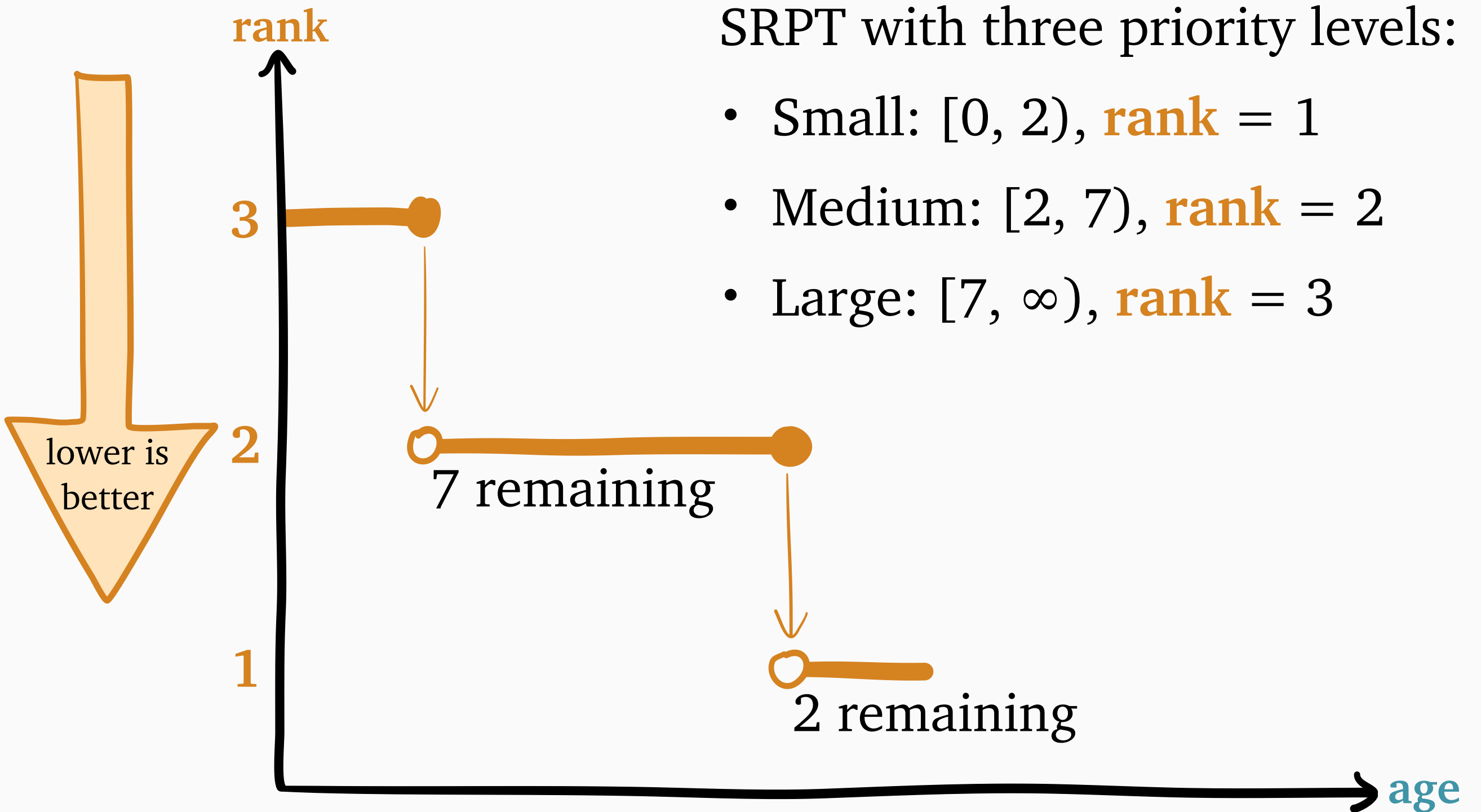
# Limited priority levels

SRPT with three priority levels:

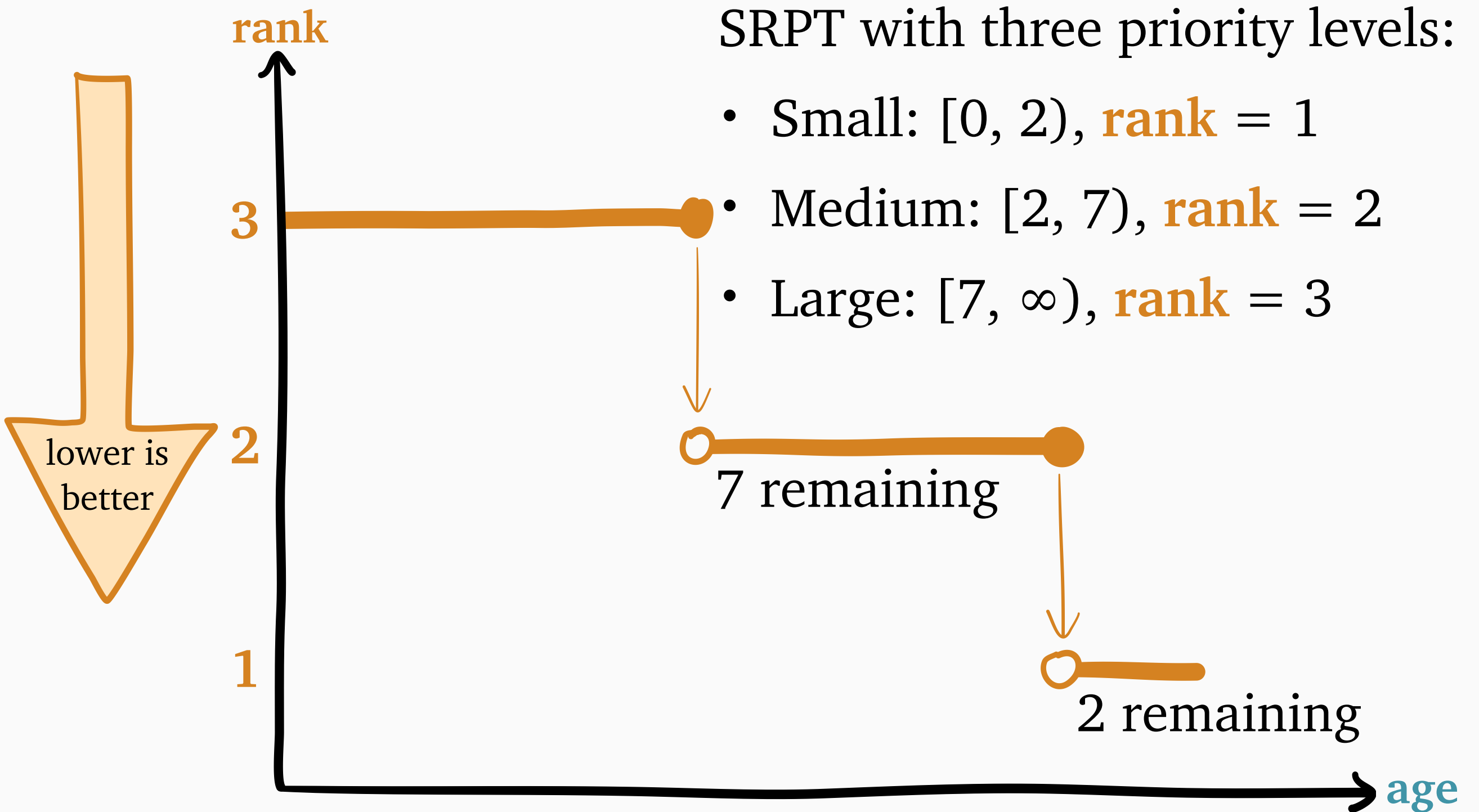
- Small:  $[0, 2)$ , **rank** = 1
- Medium:  $[2, 7)$ , **rank** = 2
- Large:  $[7, \infty)$ , **rank** = 3



# Limited priority levels

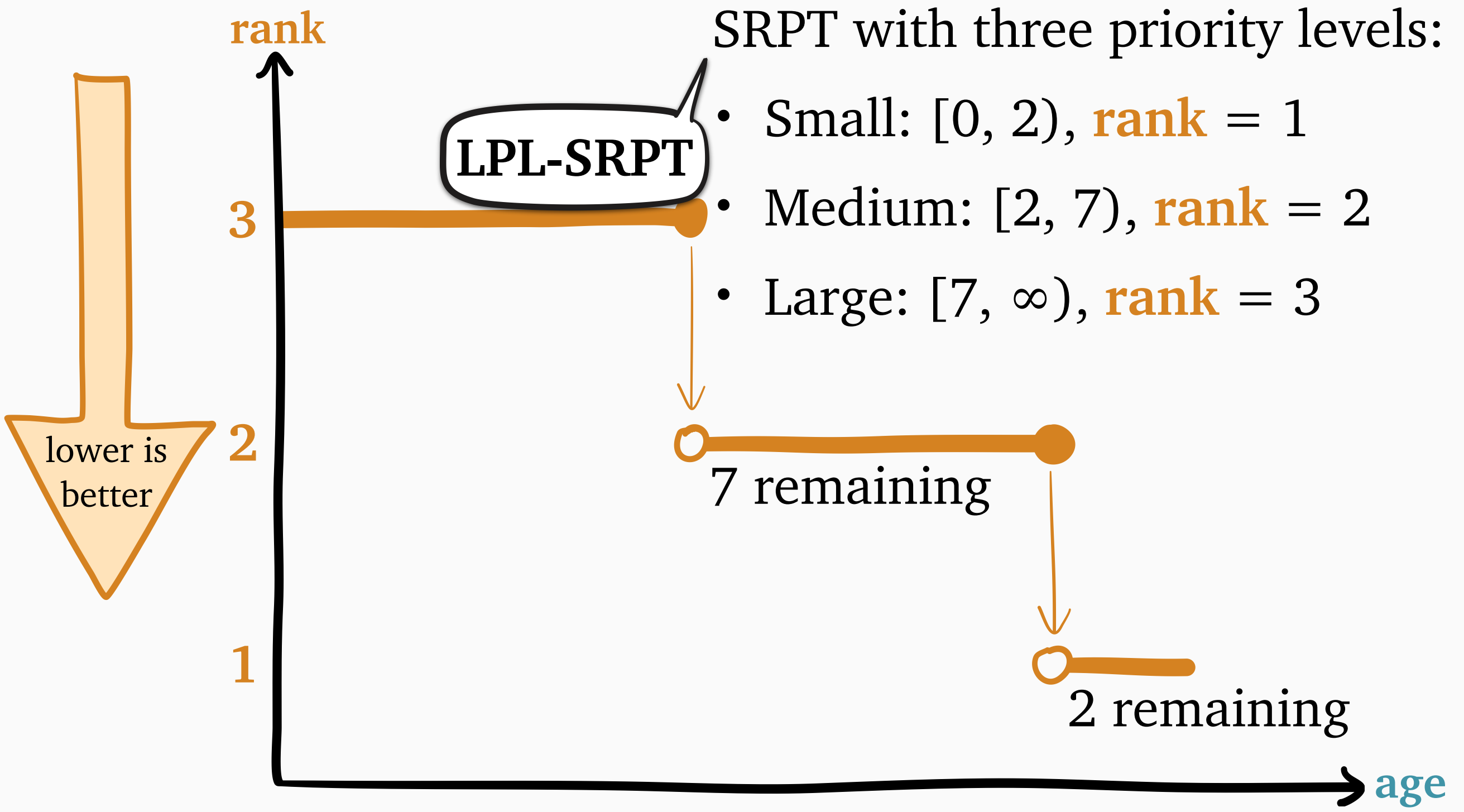


# Limited priority levels

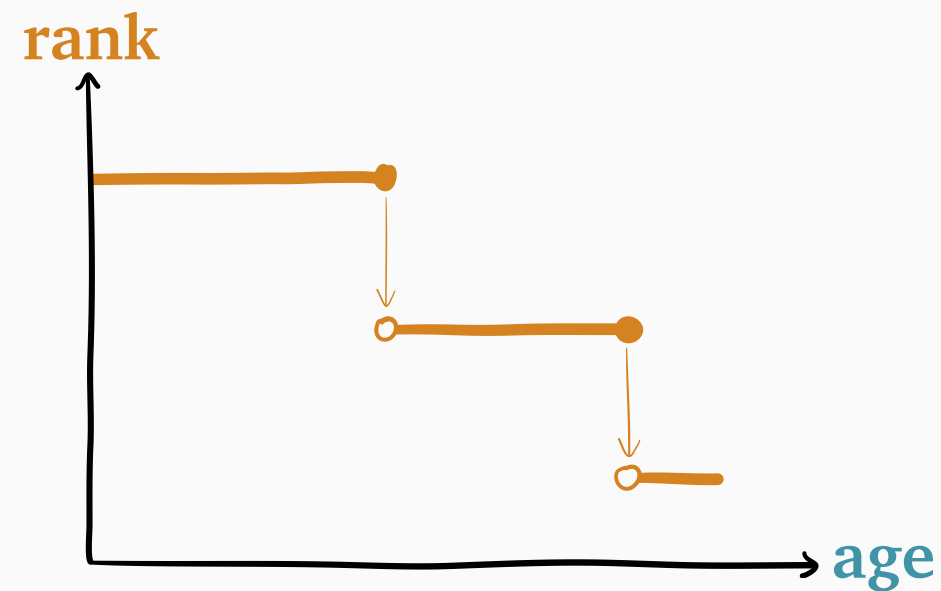




# Limited priority levels

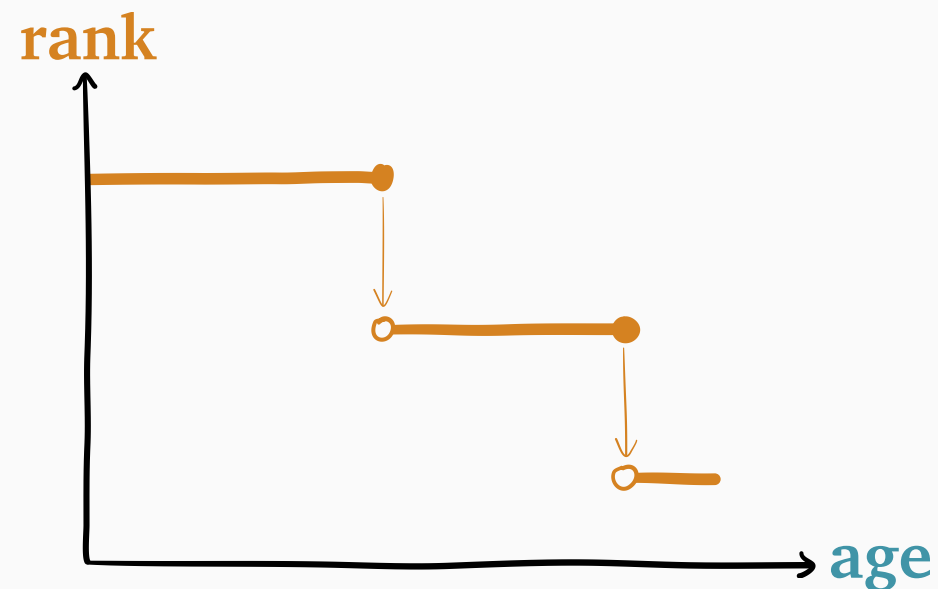


# LPL-SRPT questions



- How many levels do we need?
- How do we choose size cutoffs?
- Can we do better than LPL-SRPT?

# LPL-SRPT questions

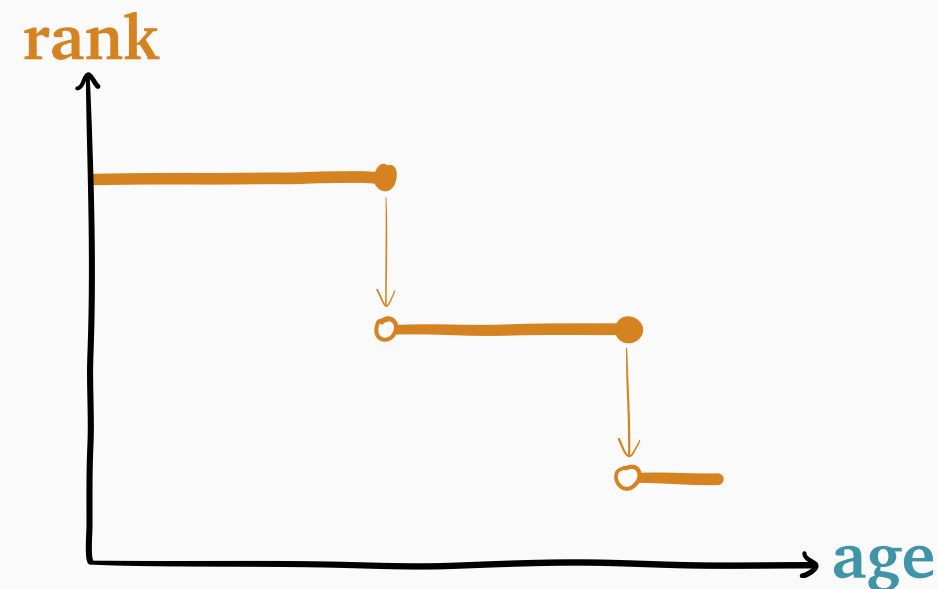


*High var: 5-ish*

*Low var: 2-ish*

- How many levels do we need?
- How do we choose size cutoffs?
- Can we do better than LPL-SRPT?

# LPL-SRPT questions



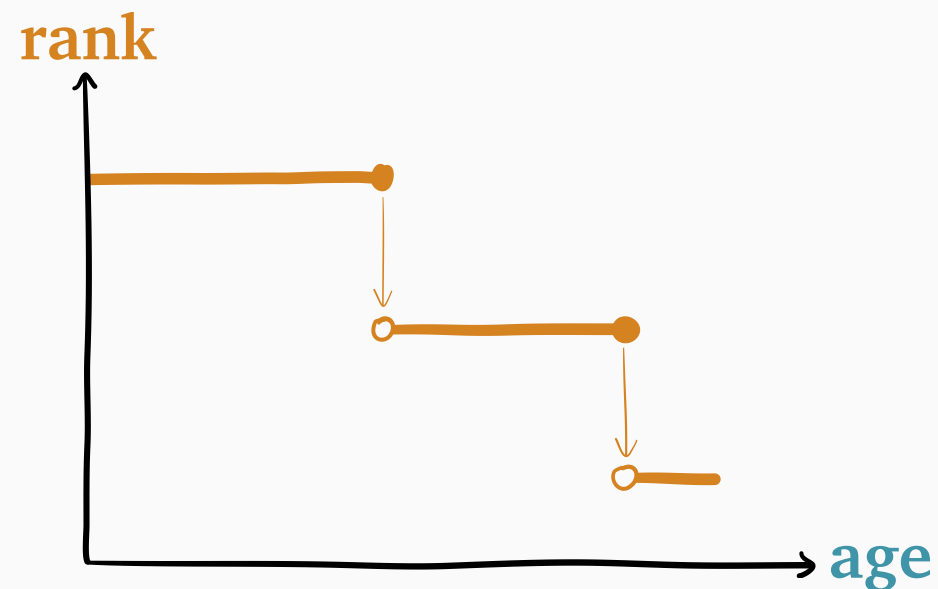
*High var: 5-ish*

*Low var: 2-ish*

Balancing load is a good heuristic

- How many levels do we need?
- How do we choose size cutoffs?
- Can we do better than LPL-SRPT?

# LPL-SRPT questions



*High var: 5-ish*

*Low var: 2-ish*

Balancing load is a good heuristic

- How many levels do we need?
- How do we choose size cutoffs?
- Can we do better than LPL-SRPT?

Yes! Constant rank better

back

*Heavy-tailed*



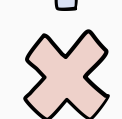
*Light-tailed*

---

*Prior work*

---

*New results*

-  = optimal
-  = intermediate
-  = pessimal

*Heavy-tailed*

*Light-tailed*

*Prior work*

✗ FCFS

✓ SRPT

✓ LAS

*New results*

✓ = optimal

? = intermediate

✗ = pessimal

*Heavy-tailed*

*Light-tailed*

*Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

*New results*

✓ = optimal

? = intermediate

✗ = pessimal



*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ SRPT

✓ LAS

✓ FCFS

✗ SRPT

✗ LAS

*New results*✓ LAS with checkpoints,  
constant gap

✓ = optimal







? = intermediate

✗ = pessimal


*Heavy-tailed*

*Light-tailed*




---

<i>Prior work</i>	 FCFS	 FCFS
	 SRPT	 SRPT
	 LAS	 LAS

---

*New results*  LAS with checkpoints,  
constant gap

not always optimal  
with growing gap

-  = optimal
-  = intermediate
-  = pessimal

*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

not always optimal  
with growing gap

*New results*

✓ LAS with checkpoints,  
constant gap

✓ **Gittins**, (**M-**)**SERPT**

✓ = optimal

? = intermediate

✗ = pessimal

*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

*New results*✓ LAS with checkpoints,  
constant gap? ✗ Non-FCFS **SOAP**✓ **Gittins, (M-)SERPT**

✓ = optimal

? = intermediate

✗ = pessimal

*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

*New results*✓ LAS with checkpoints,  
constant gap

? ✗ Non-FCFS SOAP

✓ Gittins, (M-)SERPT

✓ ? ✗ Gittins, (M-)SERPT

✓ = optimal

? = intermediate

✗ = pessimal

*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

*New results*✓ LAS with checkpoints,  
constant gap

? ✗ Non-FCFS SOAP

✓ Gittins, (M-)SERPT

✓ ? ✗ Gittins, (M-)SERPT

for some  $S$ ,  
Gittins = FCFS

✓ = optimal

? = intermediate

✗ = pessimal

*Heavy-tailed*

*Light-tailed*

*Prior work*

✗ FCFS

✓ SRPT

✓ LAS

✓ FCFS

✗ SRPT

✗ LAS

*New results*

✓ LAS with checkpoints,  
constant gap

? ✗ Non-FCFS SOAP

✓ Gittins, (M-)SERPT

? ✗ Gittins, (M-)SERPT

for some  $S$ ,  
Gittins = FCFS

for some  $S$ ,  
Gittins = LAS

✓ = optimal

? = intermediate

✗ = pessimal

*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

*New results*✓ LAS with checkpoints,  
constant gap

? ✗ Non-FCFS SOAP

✓ Gittins, (M-)SERPT

✓ ? ✗ Gittins, (M-)SERPT

✓ ? Modified Gittins

✓ = optimal

? = intermediate

✗ = pessimal



*Heavy-tailed*

*Light-tailed*

*Prior work*

✗ FCFS

✓ SRPT

✓ LAS

✓ FCFS

✗ SRPT

✗ LAS

*New results*

✓ LAS with checkpoints,  
constant gap

? ✗ Non-FCFS SOAP

✓ Gittins, (M-)SERPT

✓ ? ✗ Gittins, (M-)SERPT

✓ ? Modified Gittins

- ✓ = optimal
- ? = intermediate
- ✗ = pessimal

E[T] within  $1 + \epsilon$   
factor of optimal

*Heavy-tailed**Light-tailed**Prior work*

✗ FCFS

✓ FCFS

✓ SRPT

✗ SRPT

✓ LAS

✗ LAS

*New results*✓ LAS with checkpoints,  
constant gap

? ✗ Non-FCFS SOAP

✓ Gittins, (M-)SERPT

✓ ? ✗ Gittins, (M-)SERPT

✓ ? Modified Gittins

✓ = optimal

? = intermediate

✗ = pessimal

E[T] within  $1 + \varepsilon$   
factor of optimal

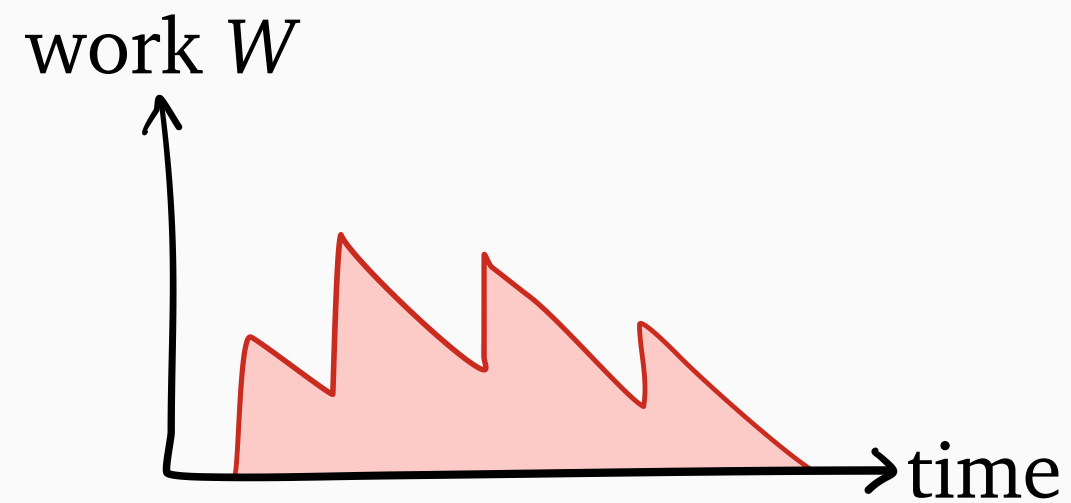
[Scully, Kreveld, Boxma, Dorsman, and Wierman, SIGMETRICS 2020]

[Scully and Kreveld, arXiv 2021]

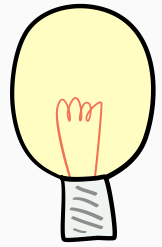
back

# Work decomposition

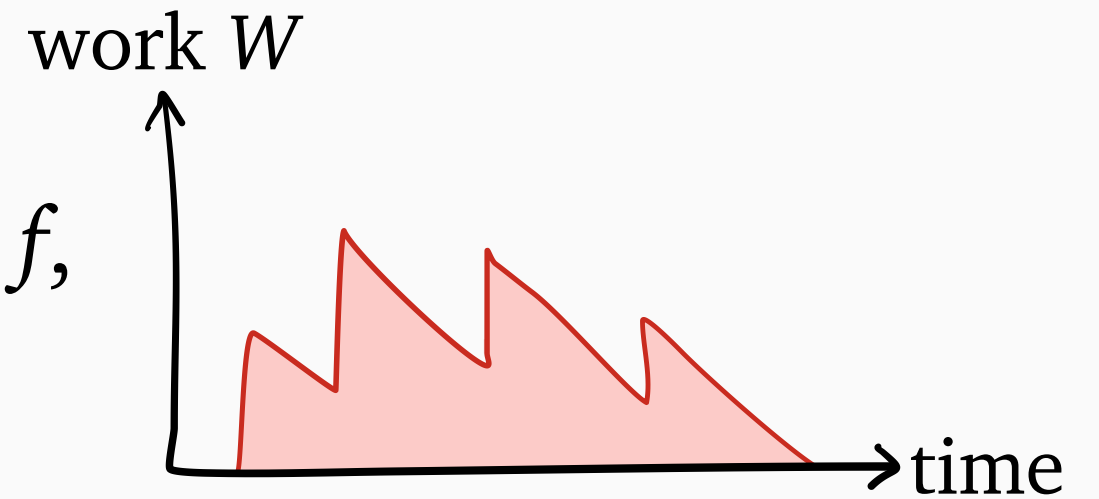
# Work decomposition



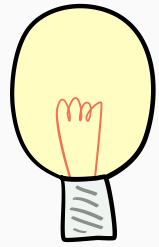
# Work decomposition



In steady-state system, for any  $f$ ,  
 $E[f(W)]$  constant w.r.t. time

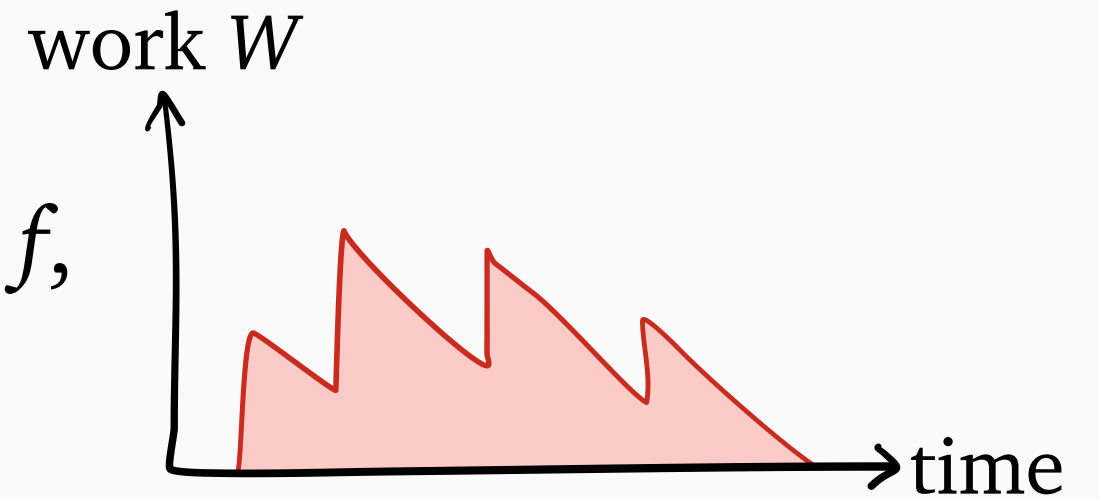


# Work decomposition



In steady-state system, for any  $f$ ,  
 $E[f(W)]$  constant w.r.t. time

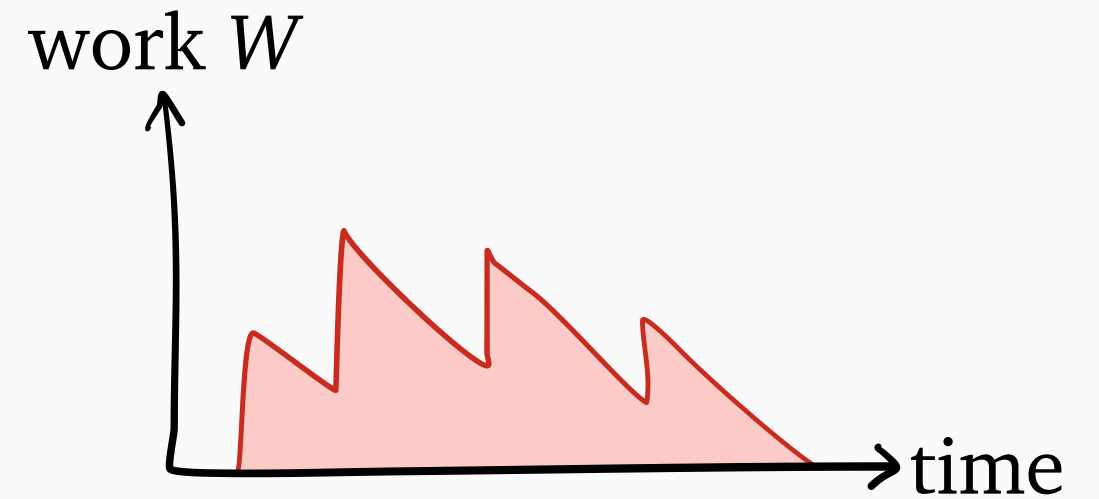
we use  
 $f(w) = w^2$



# Work decomposition

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$

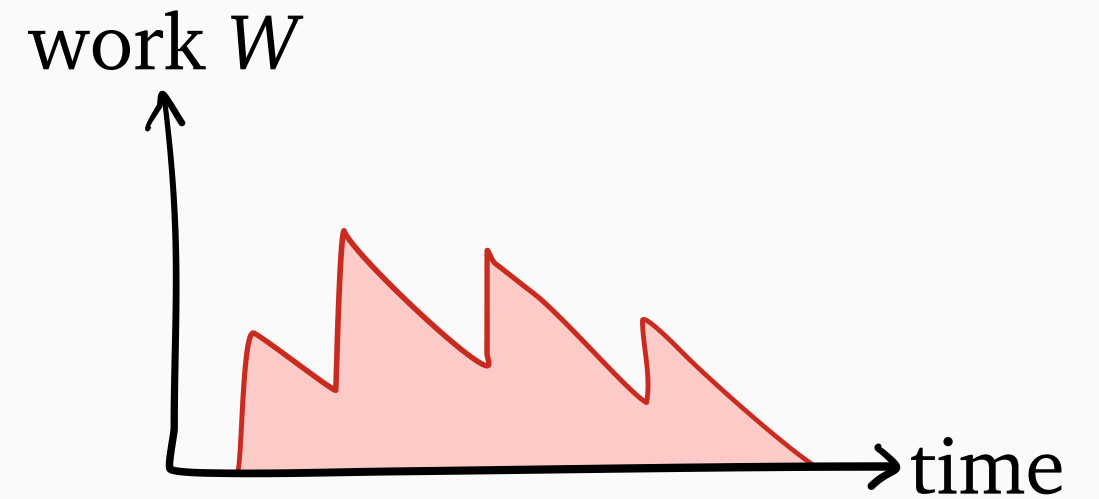


# Work decomposition

$B$  = service rate, a.k.a.  
fraction of servers busy

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$



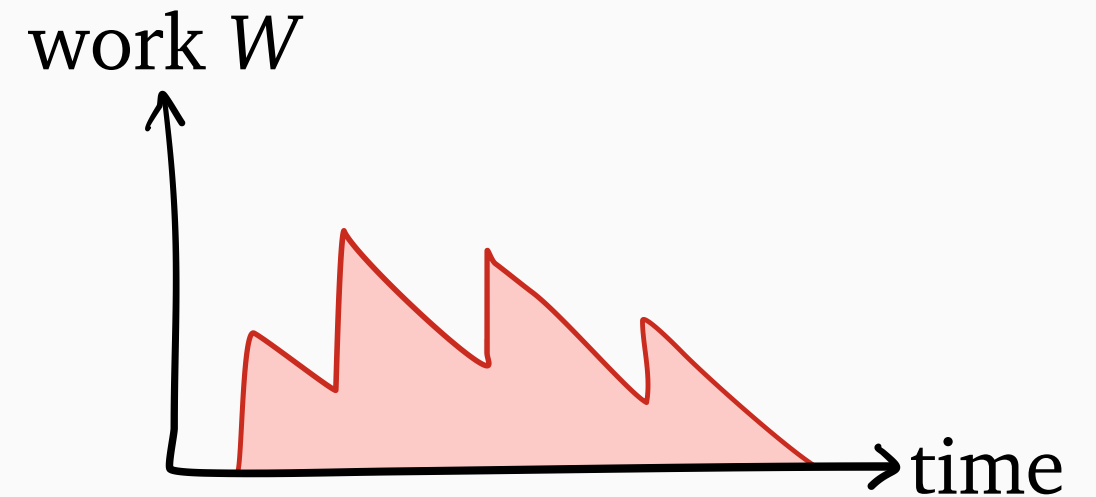


# Work decomposition

$B$  = service rate, a.k.a.  
fraction of servers busy

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$



$$\mathbf{E}[W] = \frac{\frac{\lambda}{2} \mathbf{E}[S^2]}{1 - \rho} + \frac{\mathbf{E}[(1 - B)W]}{1 - \rho}$$

# Work decomposition

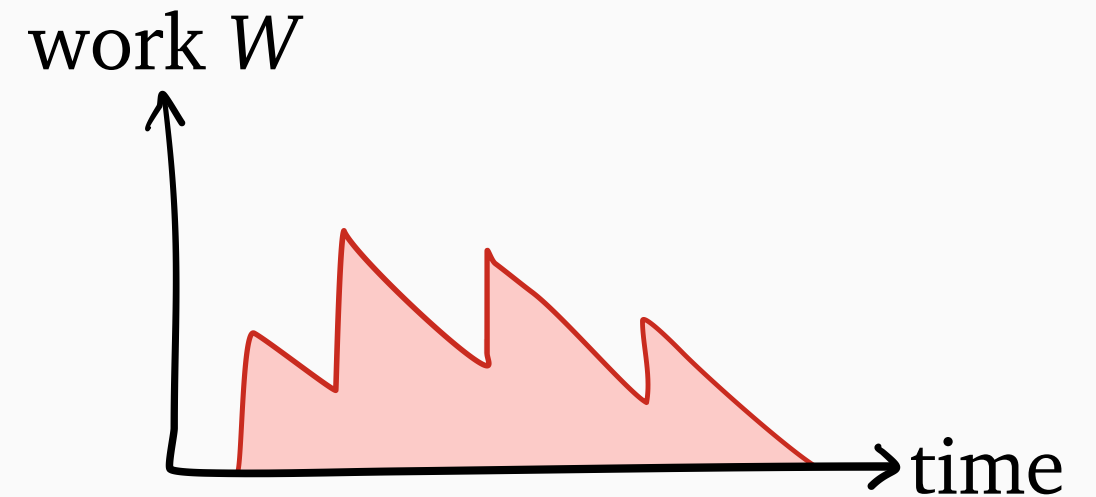
$B$  = service rate, a.k.a.  
fraction of servers busy

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$

if **single server**:  
 $(1 - B)W = 0$

$$\mathbf{E}[W] = \frac{\frac{\lambda}{2} \mathbf{E}[S^2]}{1 - \rho} + \frac{\mathbf{E}[(1 - B)W]}{1 - \rho}$$



# Work decomposition

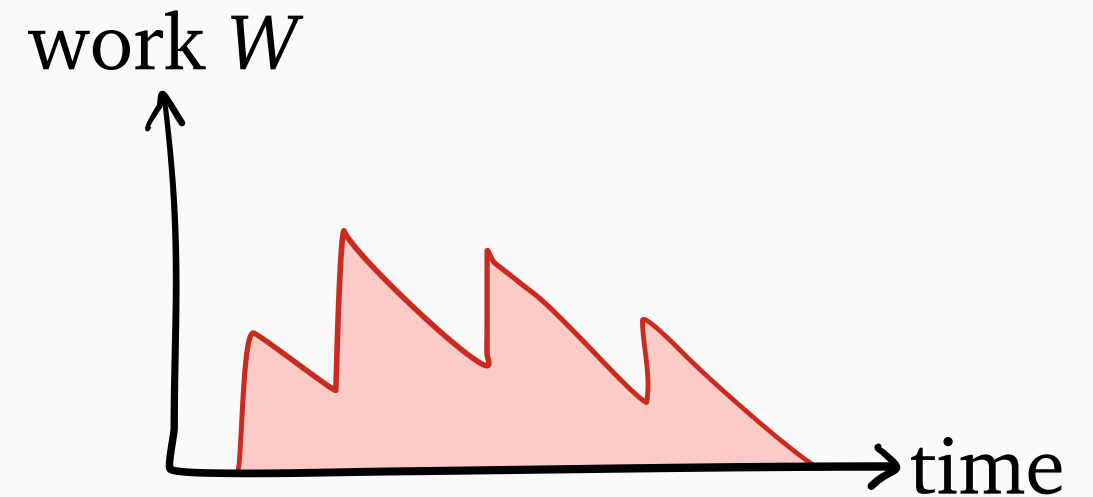
$B$  = service rate, a.k.a.  
fraction of servers busy

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$

if **single server**:  
 $(1 - B)W = 0$

$$\mathbf{E}[W] = \frac{\frac{\lambda}{2} \mathbf{E}[S^2]}{1 - \rho} + \frac{\mathbf{E}[(1 - B)W]}{1 - \rho}$$



**Lemma:**

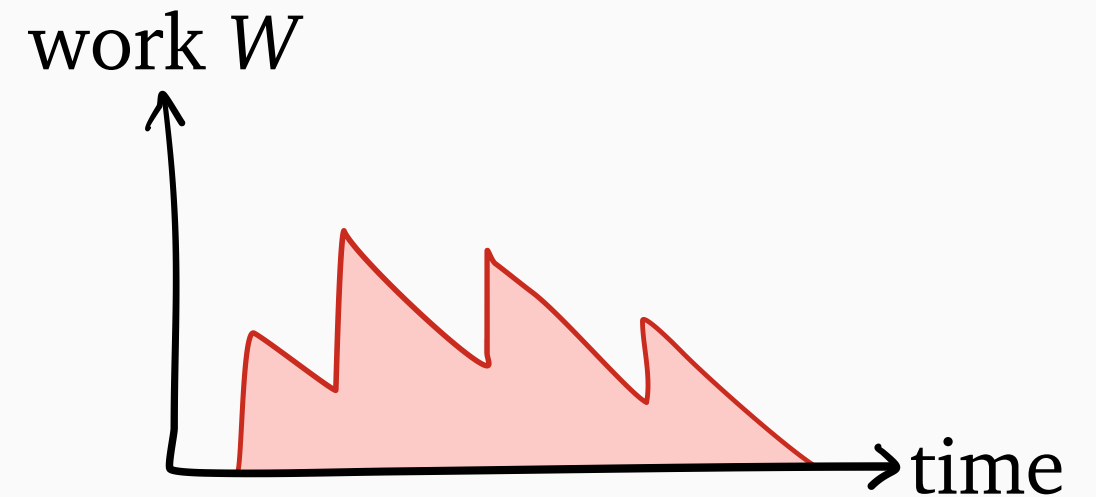
$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

# Work decomposition

$B$  = service rate, a.k.a.  
fraction of servers busy

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$



if **single server**:  
 $(1 - B)W = 0$

*newish*

$$\mathbf{E}[W] = \frac{\frac{\lambda}{2} \mathbf{E}[S^2]}{1 - \rho} + \frac{\mathbf{E}[(1 - B)W]}{1 - \rho}$$

**Lemma:**

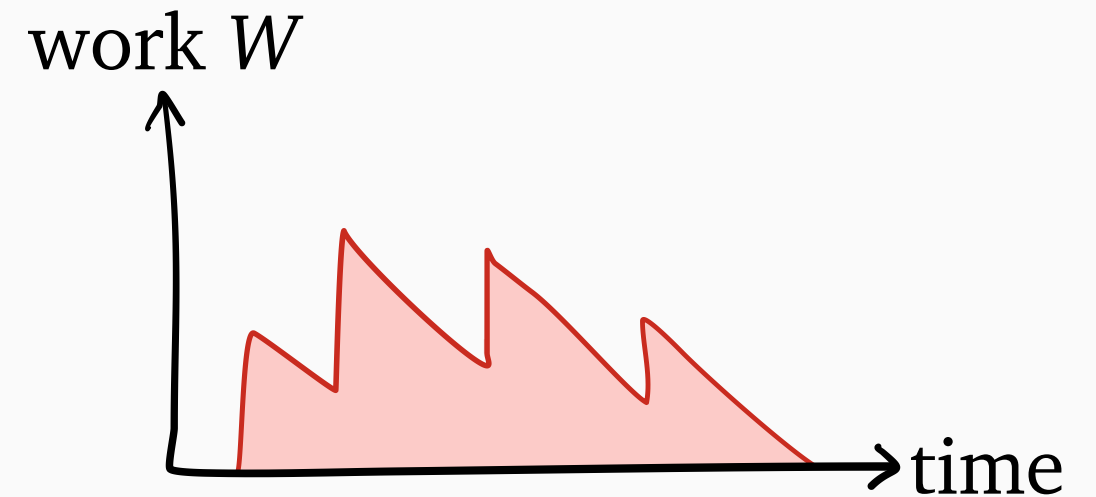
$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

# Work decomposition

$B$  = service rate, a.k.a.  
fraction of servers busy

$$\mathbf{E}[W^2 \text{ decrease rate}] = 2\mathbf{E}[BW]$$

$$\mathbf{E}[W^2 \text{ increase rate}] = \lambda \mathbf{E}[(W + S)^2 - W^2]$$



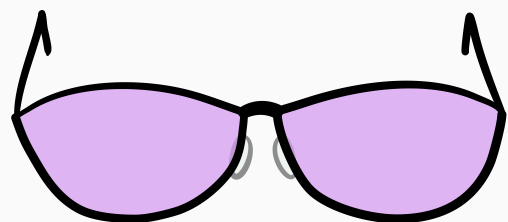
if **single server**:  
 $(1 - B)W = 0$

**newish**

$$\mathbf{E}[W] = \frac{\frac{\lambda}{2} \mathbf{E}[S^2]}{1 - \rho} + \frac{\mathbf{E}[(1 - B)W]}{1 - \rho}$$

**Lemma:**

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$



Similar story with  $r$ -work

**NEW!**

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

Suppose  $S \leq s_{\max}$  with probability 1

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

Suppose  $S \leq s_{\max}$  with probability 1

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

$\leq (k - 1)s_{\max}$



Suppose  $S \leq s_{\max}$  with probability 1

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

$\mathbf{E}[B] = \rho$

$\leq (k - 1)s_{\max}$

Suppose  $S \leq s_{\max}$  with probability 1

$$\begin{aligned} \mathbf{E}[W_k] &= \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho} \\ &\leq \mathbf{E}[W_1] + (k - 1)s_{\max} \end{aligned}$$

$\mathbf{E}[B] = \rho$

$\leq (k - 1)s_{\max}$

Suppose  $S \leq s_{\max}$  with probability 1

$$\begin{aligned} \mathbf{E}[W_k] &= \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho} \\ &\leq \mathbf{E}[W_1] + (k - 1)s_{\max} \end{aligned}$$

$\mathbf{E}[B] = \rho$

$\leq (k - 1)s_{\max}$

“work of  $\leq k - 1$  jobs”

Suppose  $S \leq s_{\max}$  with probability 1

$$\begin{aligned} \mathbf{E}[W_k] &= \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho} \\ &\leq \mathbf{E}[W_1] + (k - 1)s_{\max} \end{aligned}$$

$\mathbf{E}[B] = \rho$ 
 $\leq (k - 1)s_{\max}$

“work of  $\leq k - 1$  jobs”

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

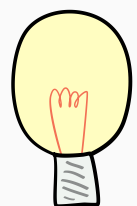
Suppose  $S \leq s_{\max}$  with probability 1

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

$$\leq \mathbf{E}[W_1] + (k - 1)s_{\max}$$

$\mathbf{E}[B] = \rho$ 
 $\leq (k - 1)s_{\max}$

“work of  $\leq k - 1$  jobs”



Single job's  $r$ -work is at most  $r$

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

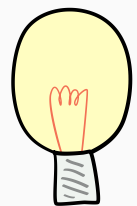
Suppose  $S \leq s_{\max}$  with probability 1

$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

$$\leq \mathbf{E}[W_1] + (k - 1)s_{\max}$$

$\mathbf{E}[B] = \rho$ 
 $\leq (k - 1)s_{\max}$

“work of  $\leq k - 1$  jobs”



Single job's  $r$ -work is at most  $r$

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

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

Suppose  $S \leq s_{\max}$  with probability 1

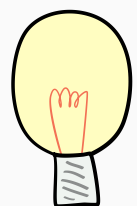
$$\mathbf{E}[W_k] = \mathbf{E}[W_1] + \frac{\mathbf{E}[(1 - B_k)W_k]}{1 - \rho}$$

$$\leq \mathbf{E}[W_1] + (k - 1)s_{\max}$$

$$\mathbf{E}[B] = \rho$$

$$\leq (k - 1)s_{\max}$$

“work of  $\leq k - 1$  jobs”



Single job's  $r$ -work is at most  $r$

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

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

can improve

# SRPT- $k$ and Gittins- $k$ $E[T]$



**WINE** gives *first analysis* of

- SRPT- $k$
- Gittins- $k$



# SRPT- $k$ and Gittins- $k$ $E[T]$



**WINE** gives *first analysis* of

- SRPT- $k$
- Gittins- $k$

**Theorem:** for SRPT and Gittins,

$$E[T_k] \leq E[T_1] + (k - 1) \cdot O\left(\log \frac{1}{1 - \rho}\right)$$

# SRPT- $k$ and Gittins- $k$ $E[T]$



**WINE** gives *first analysis* of

- SRPT- $k$
- Gittins- $k$

**Theorem:** for SRPT and Gittins,

$$E[T_k] \leq E[T_1] + \underbrace{(k-1) \cdot O\left(\log \frac{1}{1-\rho}\right)}_{o(E[T_1])}$$

# SRPT- $k$ and Gittins- $k$ $E[T]$



**WINE** gives *first analysis* of

- SRPT- $k$
- Gittins- $k$

**Theorem:** for SRPT and Gittins,

$$E[T_k] \leq E[T_1] + \underbrace{(k-1) \cdot O\left(\log \frac{1}{1-\rho}\right)}_{o(E[T_1])}$$

**Corollary:** SRPT and Gittins minimize  $E[T]$  in heavy traffic (in their respective settings)