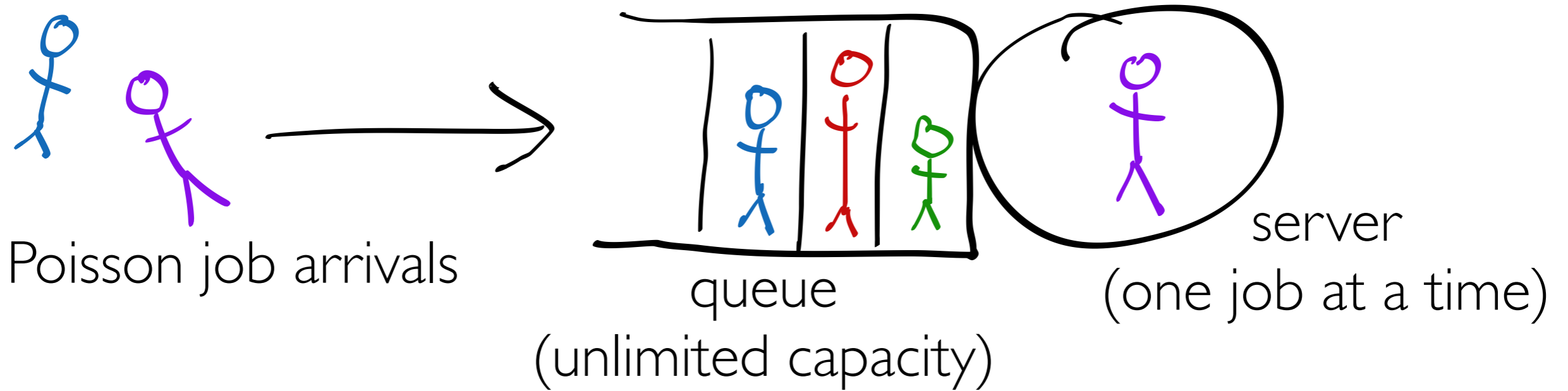


Scheduling

with the

Gittins Index

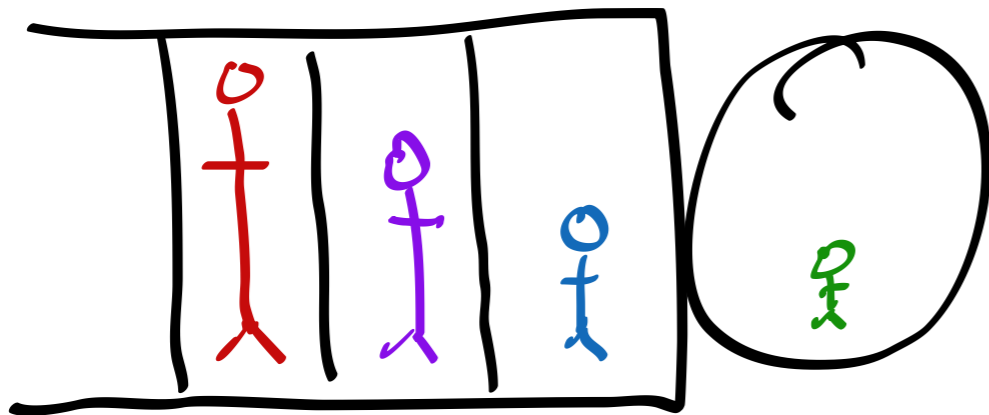
Ziv Scully
CMU Theory Lunch



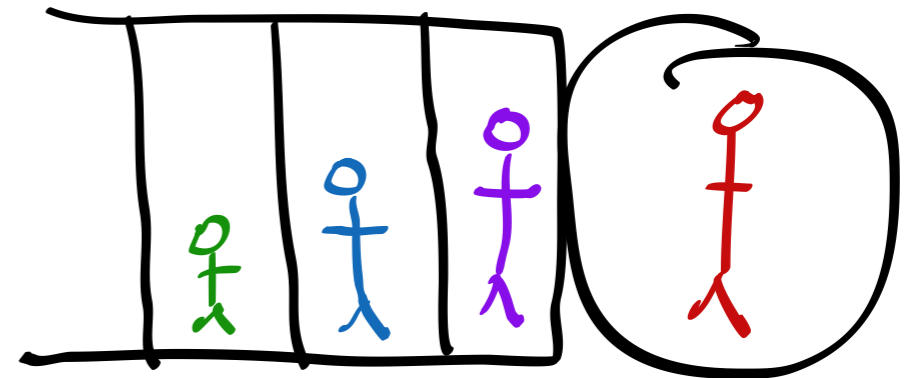
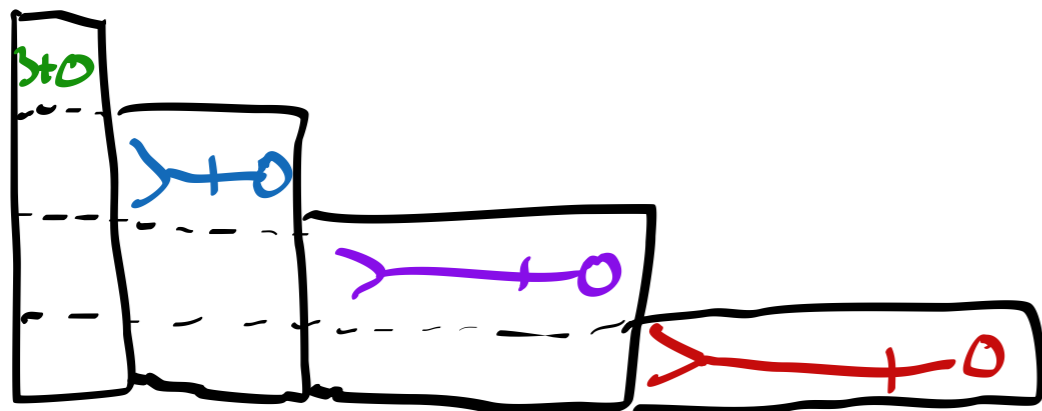
Which job should we serve to minimize mean response time?

Which job should we serve to minimize mean response time?

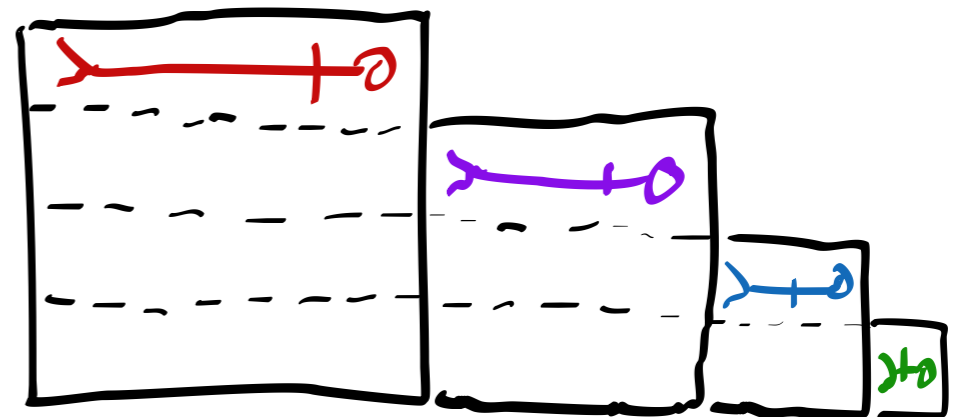
The shortest one!



Total time in system:

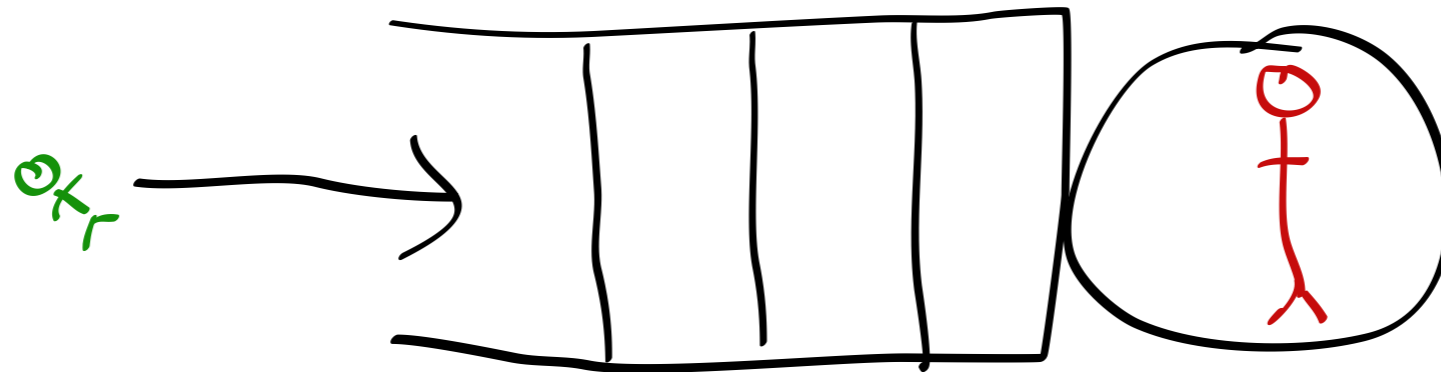


Total time in system:



But what if...

... a new job arrives?

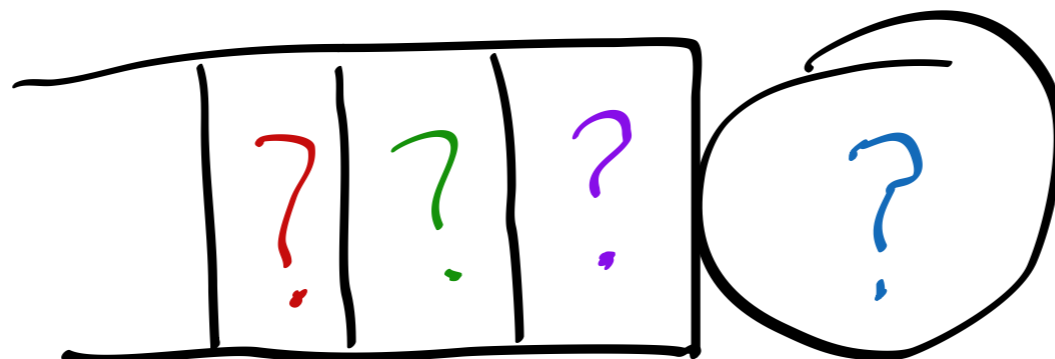


If preemptive (we can interrupt jobs):
shortest remaining processing time
(SRPT)

If nonpreemptive:
shortest job first (SJF)

But what if...

... we don't know the job sizes?



Do jobs one at a time?

FIFO, LIFO, random

Do multiple jobs at once?

processor sharing,
foreground-background

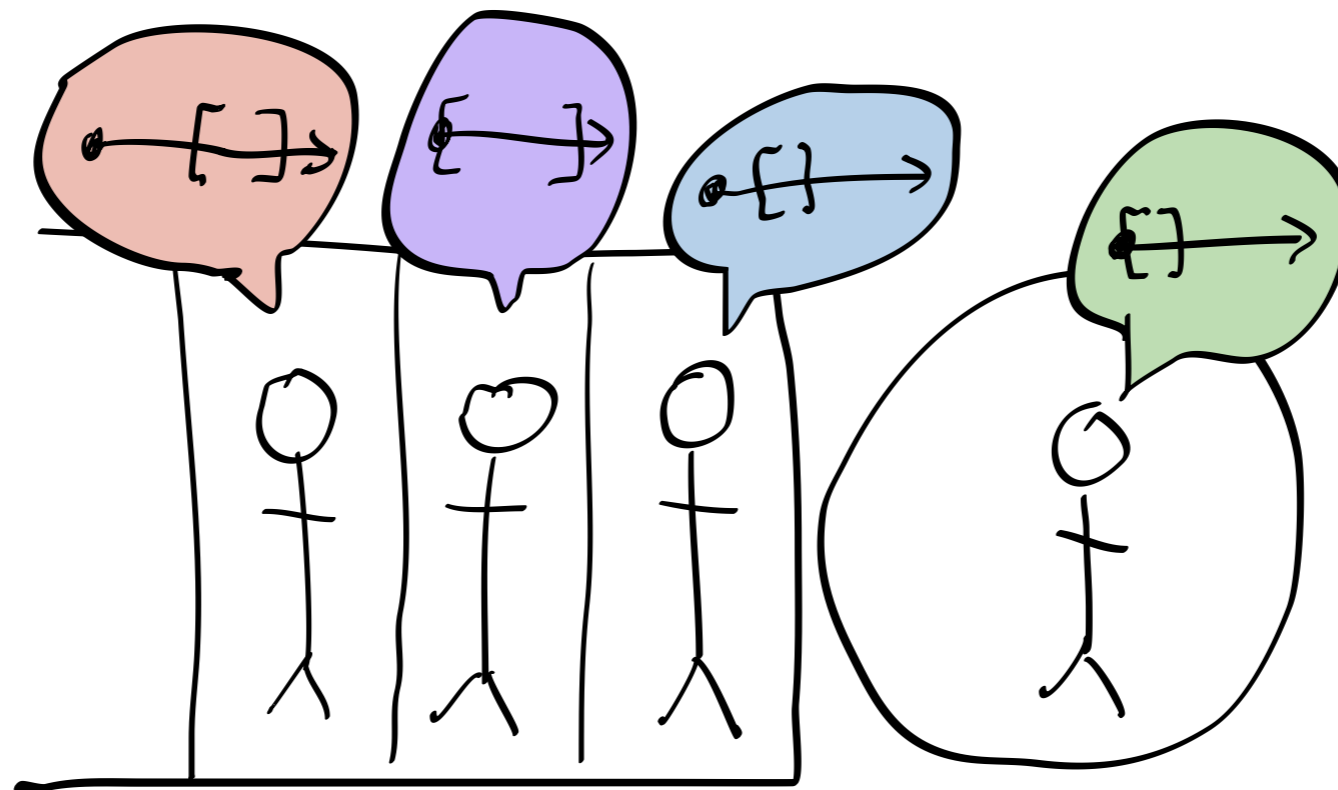
Use job size distribution?

SERPT (like SRPT, E for expected),

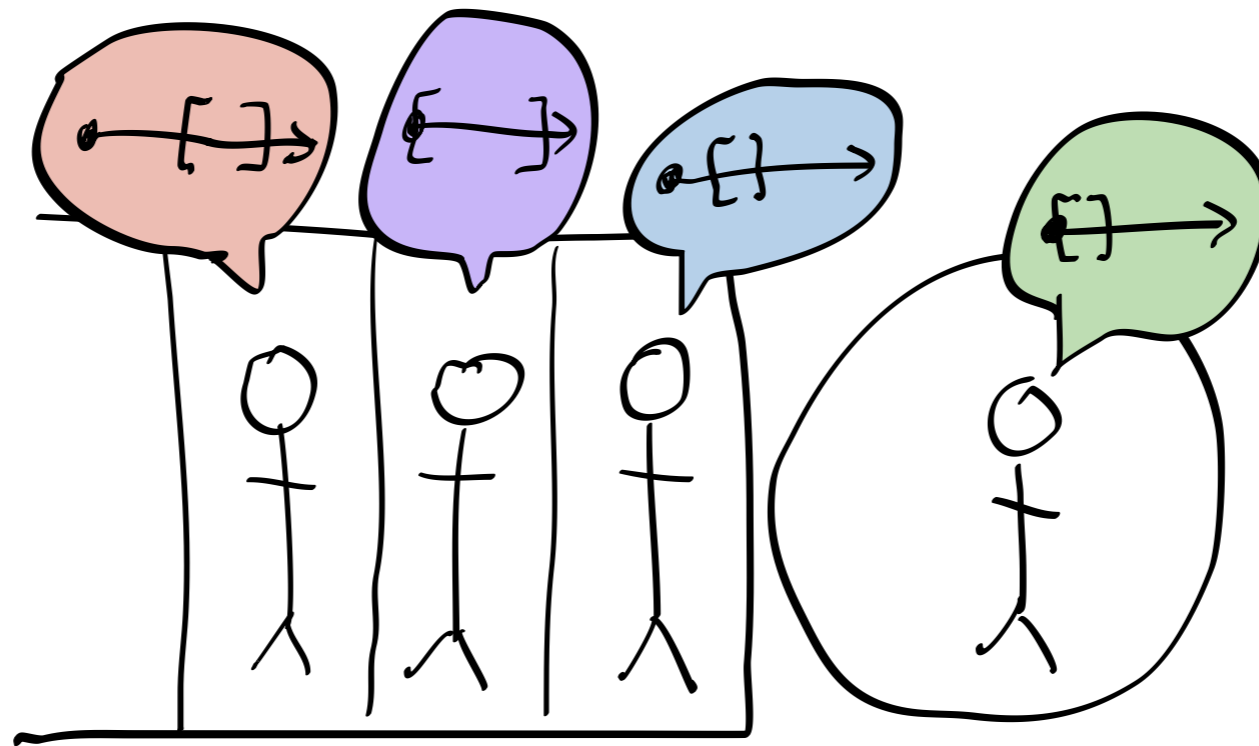
highest hazard rate

How should we schedule
in a preemptive setting
with unknown job sizes?

Example: uniform job size distributions



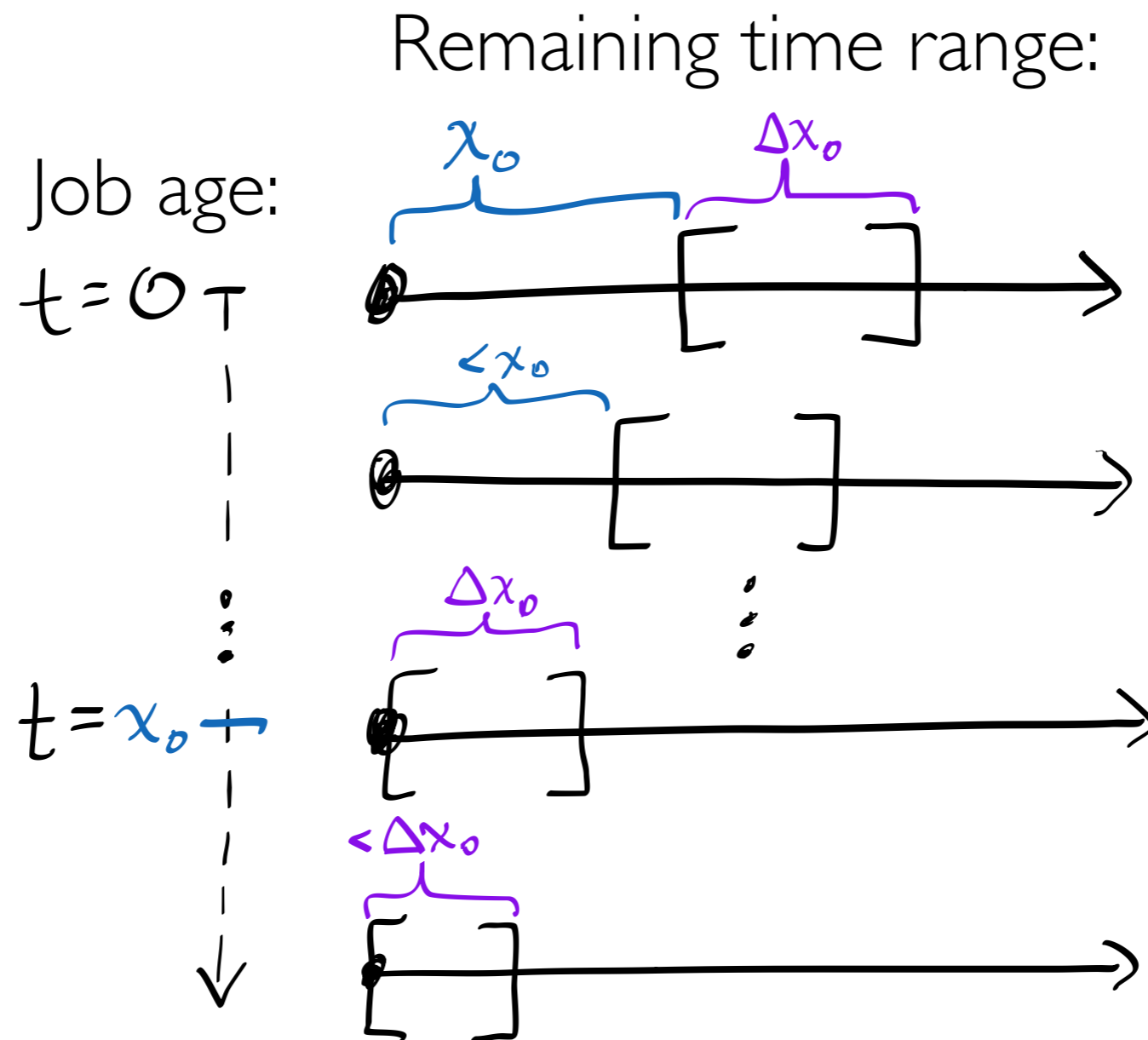
Guess: SERPT is best for uniformly distributed jobs



Other than arrivals,
when do we preempt?

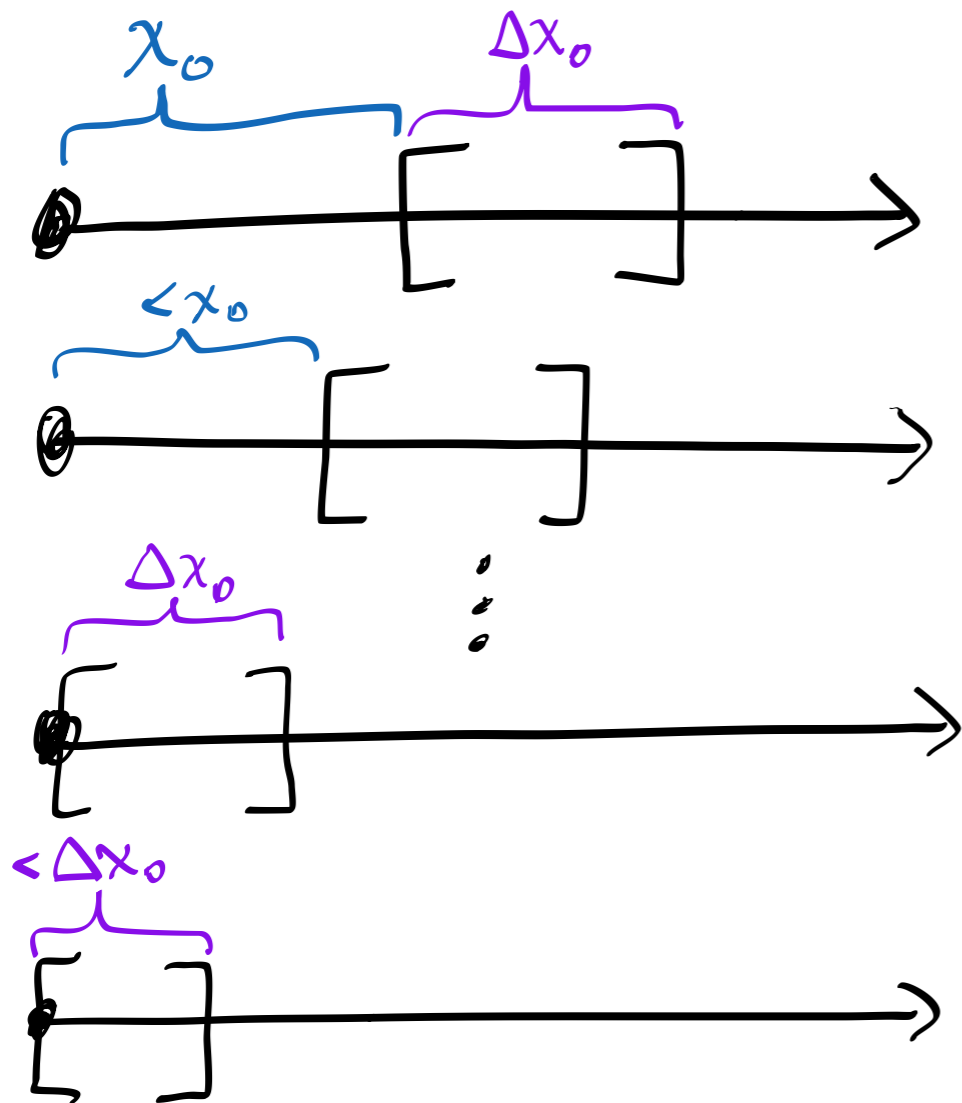
Other than arrivals,
when do we preempt?

Never!



Expected remaining time:

$$E[X(t)] = x(t) + \frac{\Delta x(t)}{2}$$



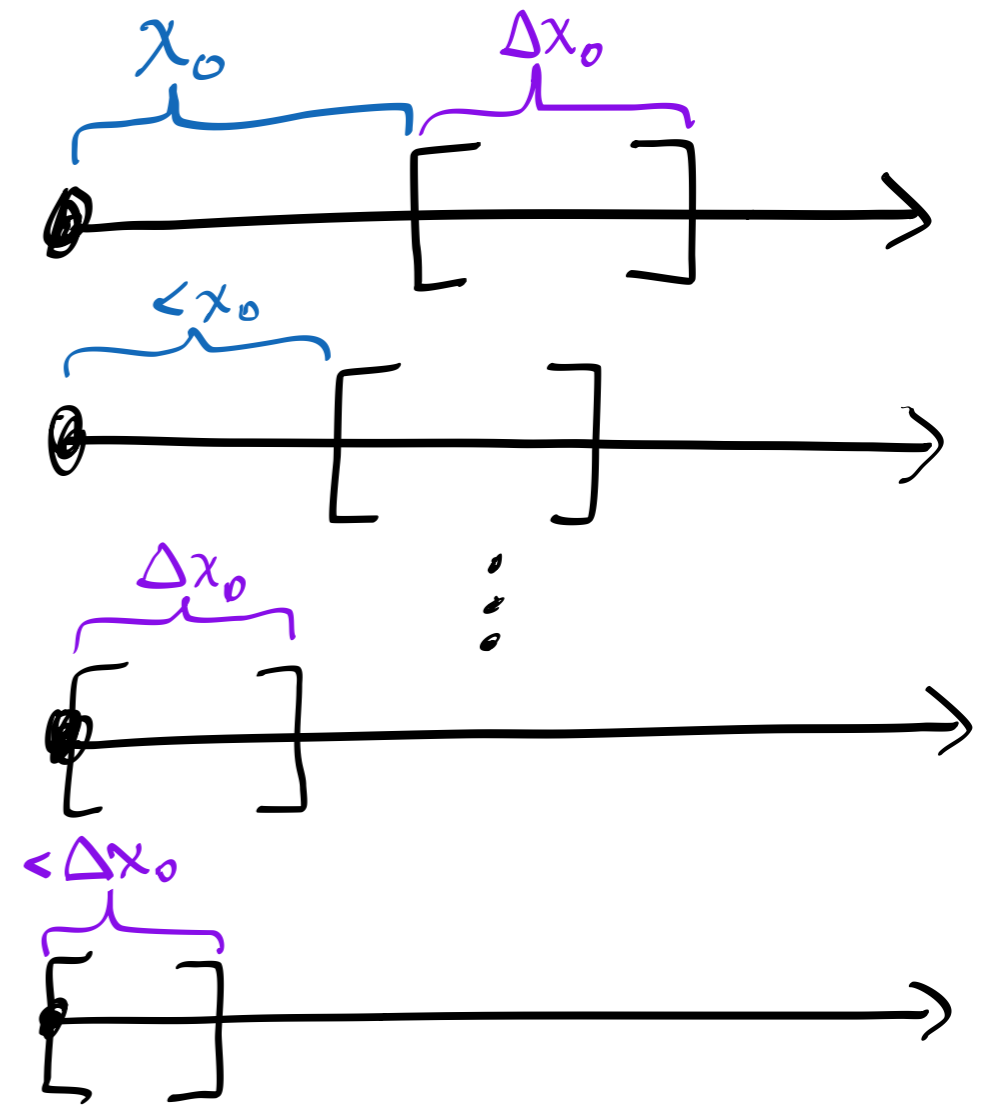
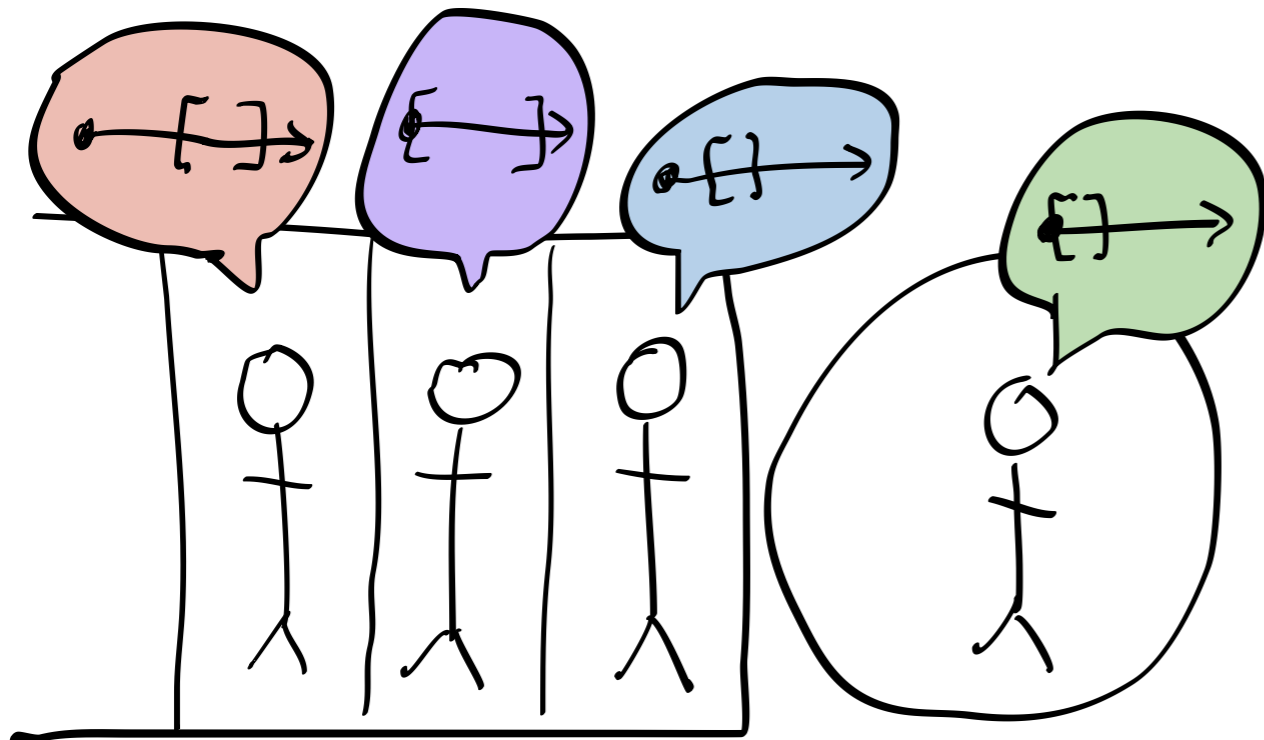
$$X(t) \sim [X - t \mid X > t]$$

$$x(t) = (x_0 - t)^+ \quad \hat{\curvearrowright}$$

$$\Delta x(t) = (\Delta x_0 - (t - x_0)^+)^+ \quad \hat{\curvearrowright}$$

Because remaining times get shorter with age...

... we never actively look to preempt jobs.

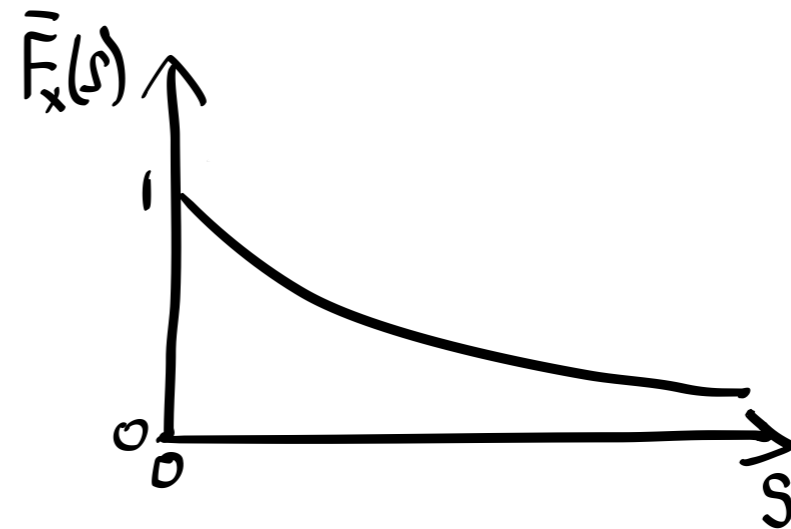
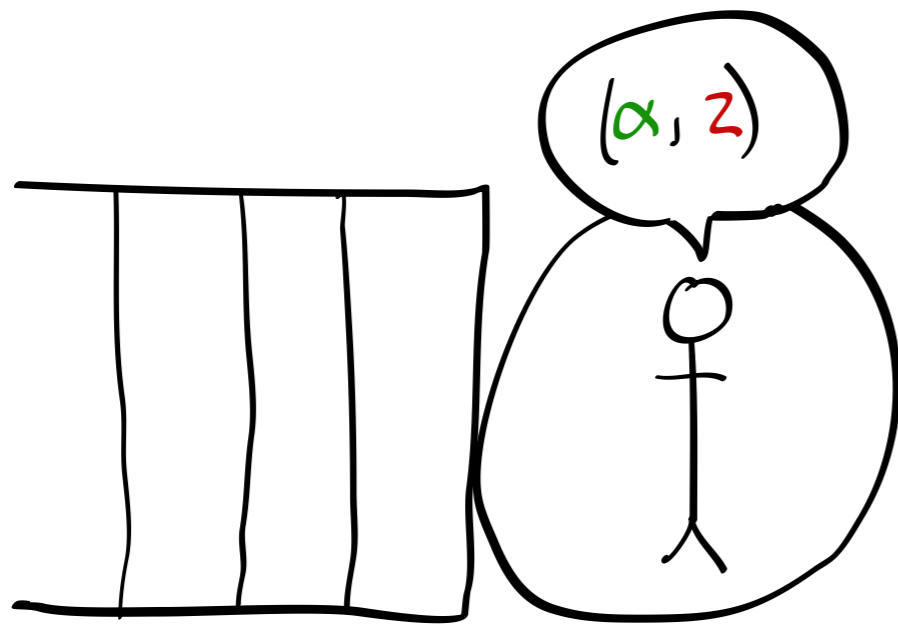


SERPT is optimal when all jobs get shorter with age (DMRL).

Example: Pareto job size distributions

Tail of size distribution:

$$\bar{F}_x(s) = \left(1 + \frac{s}{z}\right)^{-\alpha}$$



Say we do SERPT.
When do we preempt?

$$X(t) \sim [X-t | X > t]$$

$$\bar{F}_{X(t)}(s) = \frac{\bar{F}_X(s+t)}{\bar{F}_X(t)} = \left(1 + \frac{s}{z+t}\right)^{-\alpha}$$

To put that another way:

$$z(t) = z_0 + t$$

$$z \uparrow \Rightarrow 1 + \frac{s}{z} \downarrow \Rightarrow \bar{F}_X(s) \uparrow$$

\Rightarrow Jobs get longer with age!

Pareto has *decreasing hazard rate*

$$h_x(t) = \frac{f_x(t)}{\bar{F}_x(t)} = \frac{\alpha}{z+t} \quad (\text{for Pareto})$$

Hazard rate is like conditional density:
having reached some age, how likely
is it that the job finishes right now?

decreasing hazard rate \Rightarrow older jobs less likely to finish
 \Rightarrow lots of preemption

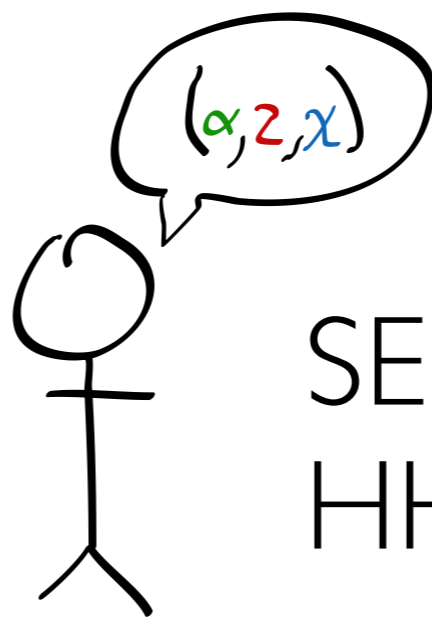
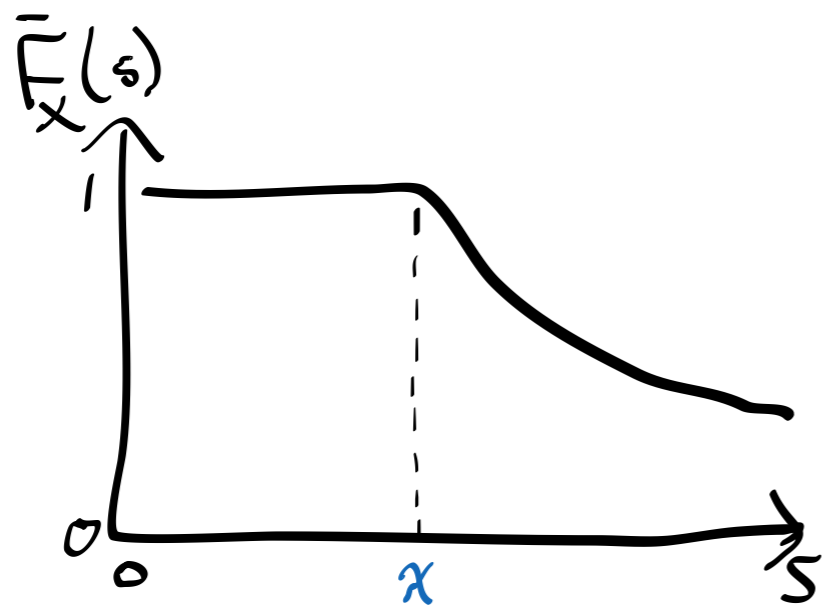
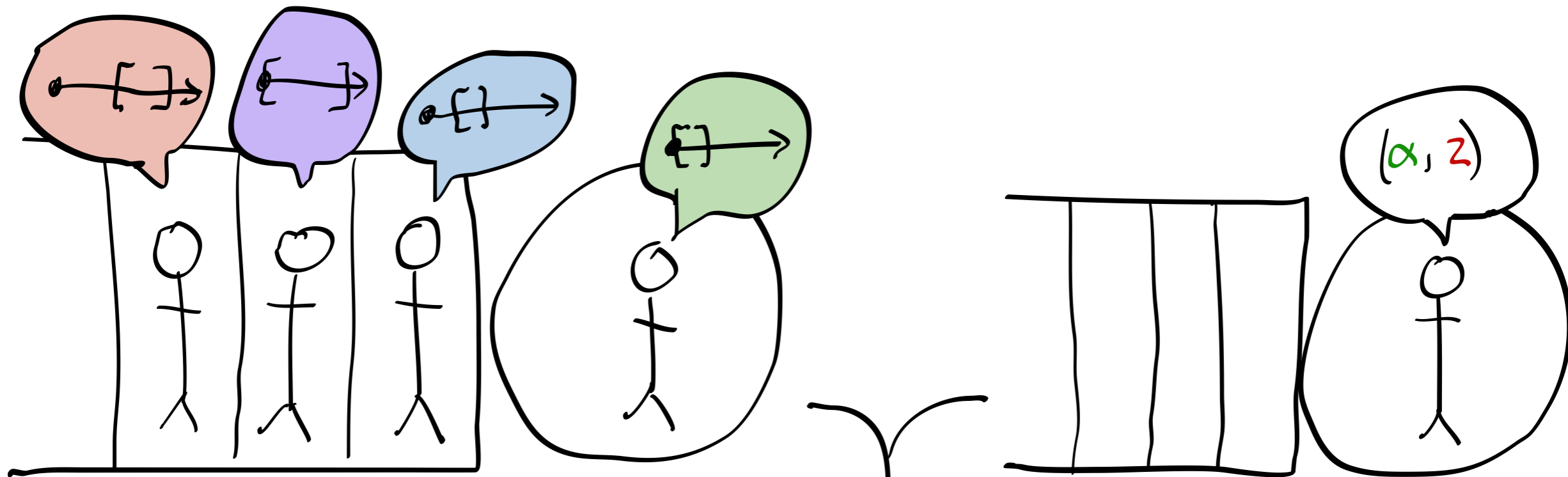
Highest hazard rate (HHR):
always serve job of highest hazard rate

fixed $\alpha \Rightarrow$ equivalent to SERPT

mixed $\alpha \Rightarrow$ handles tradeoffs differently

$$\begin{aligned} \because \alpha &> \alpha' \\ z &> z' \quad \because \end{aligned}$$

HHR is optimal when all jobs have decreasing hazard rate (DHR).



SERPT? Not DMRL.
HHR? Not DHR.

Two notions of completion rate

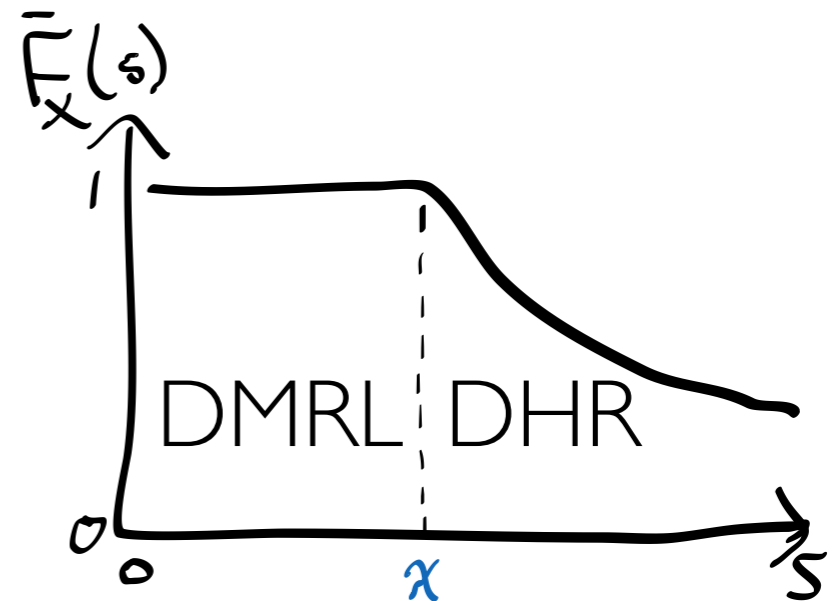
SERPT: $\frac{1}{E[X(t)]} = \frac{\int_t^\infty f_x(s) ds}{\int_t^\infty \bar{F}_x(s) ds}$

right idea if we *never* look for preemptions (like DMRL case)

HHR: $h_x(t) = \lim_{T \rightarrow t} \frac{\int_t^T f_x(s) ds}{\int_t^T \bar{F}_x(s) ds}$

right idea if we *always* look for preemptions (like DHR case)

What if we only *sometimes* look for preemptions?



Gittins index: best possible completion rate

$$v_x(t) = \sup_{T > t} \frac{P\{X(t) < T\}}{E[\min\{X(t), T\}]} = \sup_{T > t} \frac{\int_t^T f_x(s) ds}{\int_t^T \bar{F}_x(s) ds}$$

$$\frac{1}{E[X(t)]} = \frac{\int_t^\infty f_x(s) ds}{\int_t^\infty \bar{F}_x(s) ds}$$

$$h_x(t) = \lim_{T \rightarrow t} \frac{\int_t^T f_x(s) ds}{\int_t^T \bar{F}_x(s) ds}$$