

Stochastic Scheduling with Predictions

Isaac Grosf

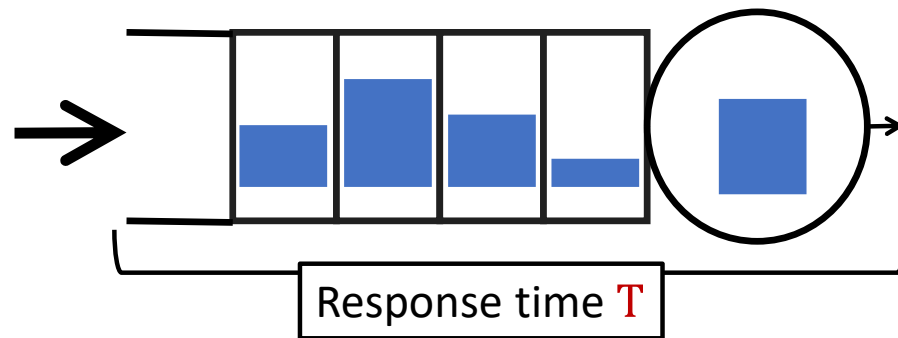
CMU Computer Science

Ziv Scully (CMU -> Cornell)

Michael Mitzenmacher (Harvard)

Scheduling with Predictions

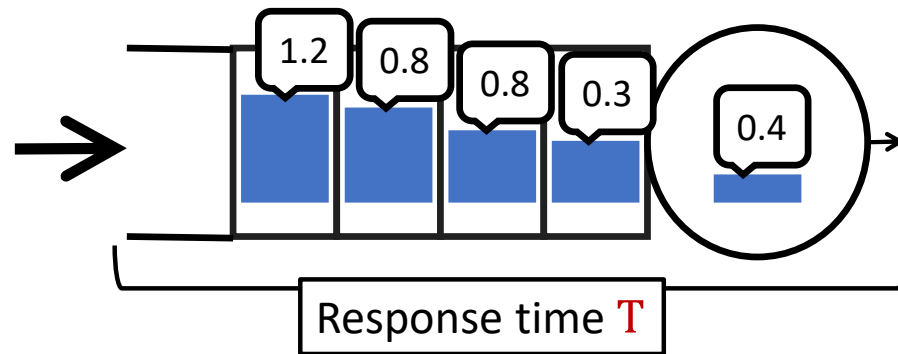
“Uniform Bounds for Scheduling with Job Size Estimates” ITCS 2022



Goal: Minimize mean response time ($E[T]$)

Scheduling with Predictions

“Uniform Bounds for Scheduling with Job Size Estimates” ITCS 2022



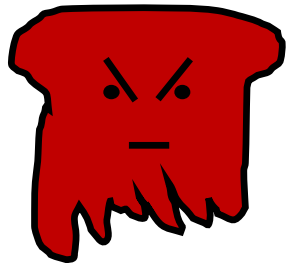
Goal: Minimize mean response time ($E[T]$)

Known sizes: Shortest Remaining Processing Time (SRPT) is optimal

Predicted sizes: ?

Stochastic Scheduling with Predictions

Two ways to study:



Worst Case Analysis



Stochastic Analysis

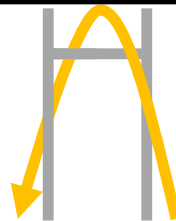
This talk

Why stochastic analysis for scheduling?

Reflects real world



Overcomes worst-case barriers

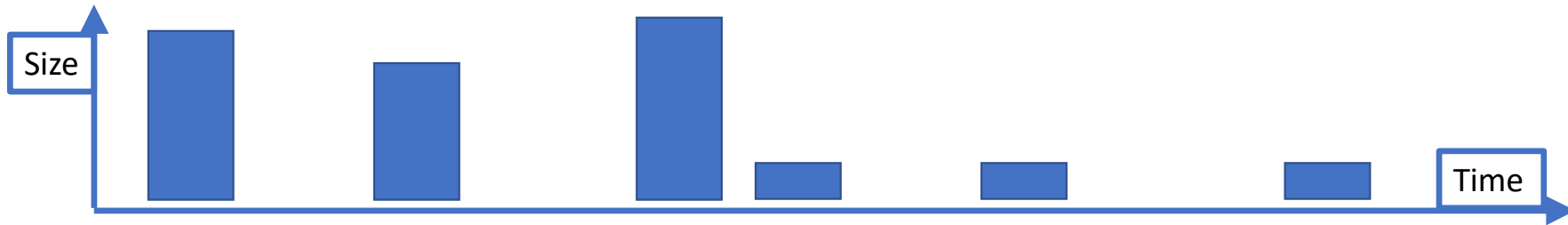


Challenging, subtle problems



Why Stochastic Arrivals?

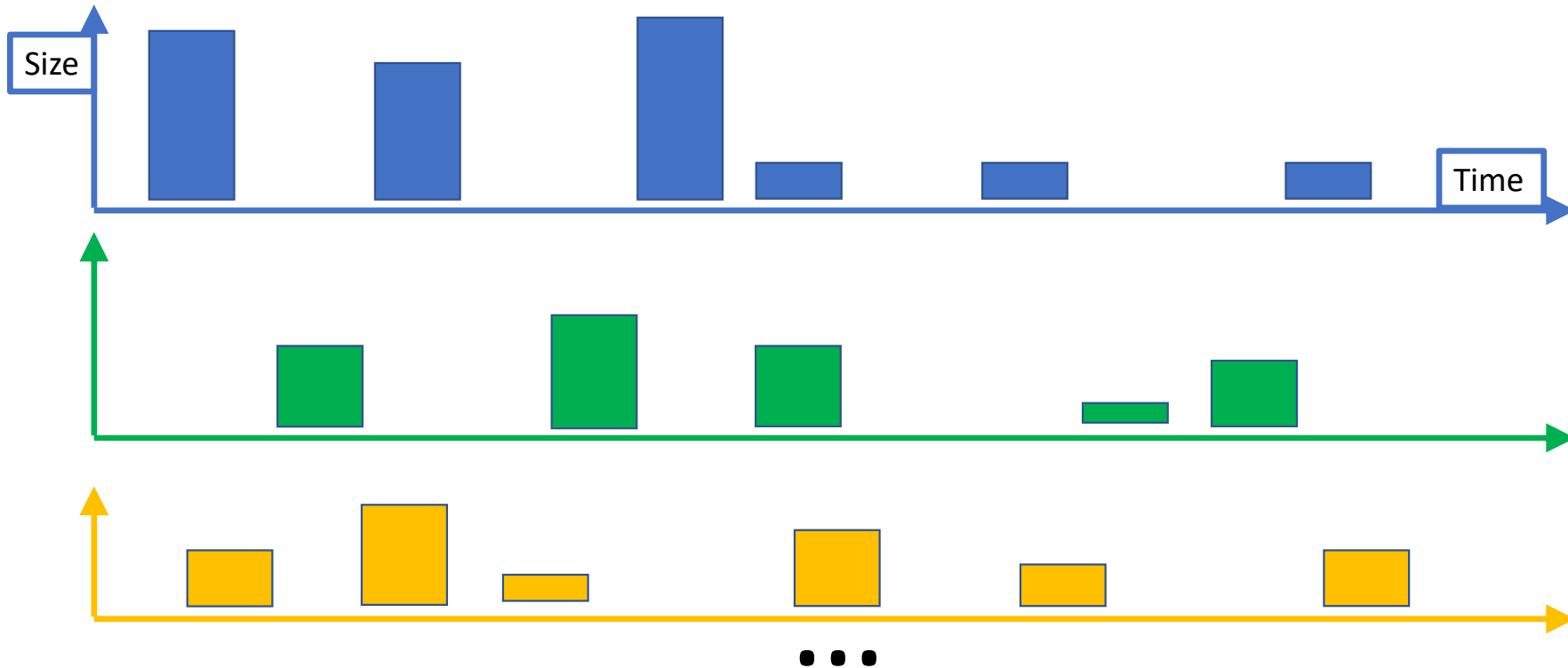
Worst case setting: General arrival sequence



Why Stochastic Arrivals?

Worst case setting: General arrival sequence

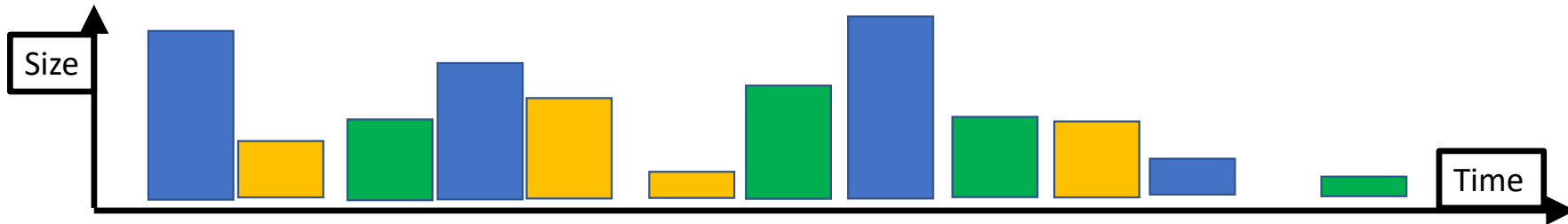
 Real world: lots of independent arrival sequences



Why Stochastic Arrivals?

Worst case setting: General arrival sequence

 Real world: lots of independent arrival sequences

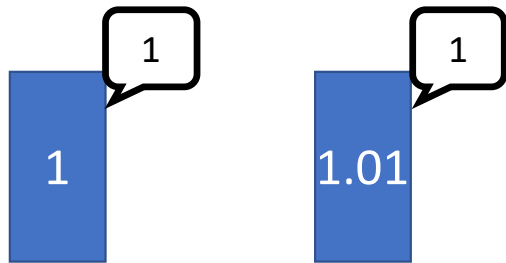
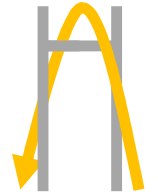


Approximation for large systems:

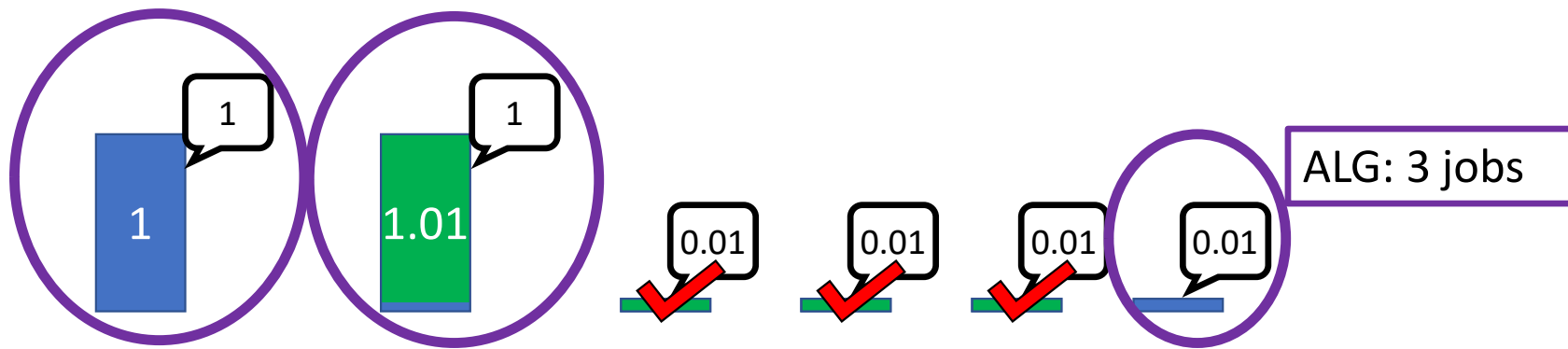
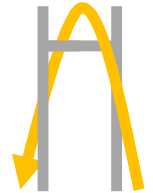
Exponential interarrival times (Poisson)

I.i.d. sizes

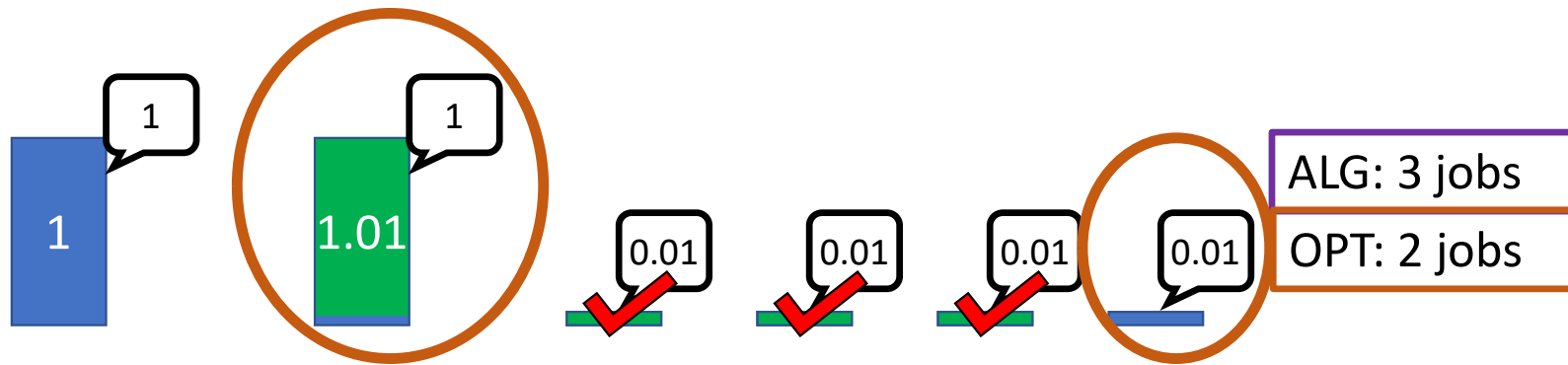
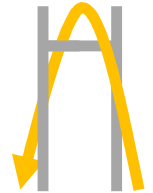
Example of Scheduling with Predictions



Example of Scheduling with Predictions



Example of Scheduling with Predictions

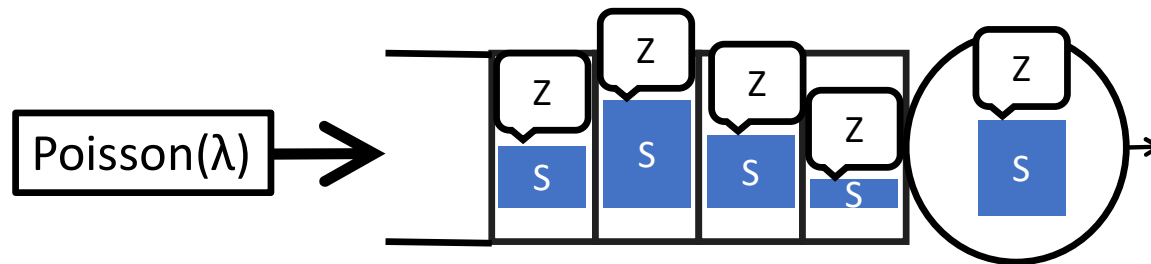


Tiny prediction error \rightarrow $3/2$ competitive ratio

Very unlikely in our stochastic model

Goal: Tiny prediction error \rightarrow Near-perfect performance

Specific model



Arrival process: Poisson(λ)

Size and predicted size i.i.d. from (S, Z)

$Z \in [\beta S, \alpha S]$

(S, Z) distribution chosen adversarially

Scheduler does not know S, Z, α, β

Metric: $\frac{E[T^\pi]}{E[T^{OPT}]} = \frac{E[T^\pi]}{E[T^{SRPT}]}$

Performance Goals

Goals for $\frac{E[T^\pi]}{E[TSRPT]}$:

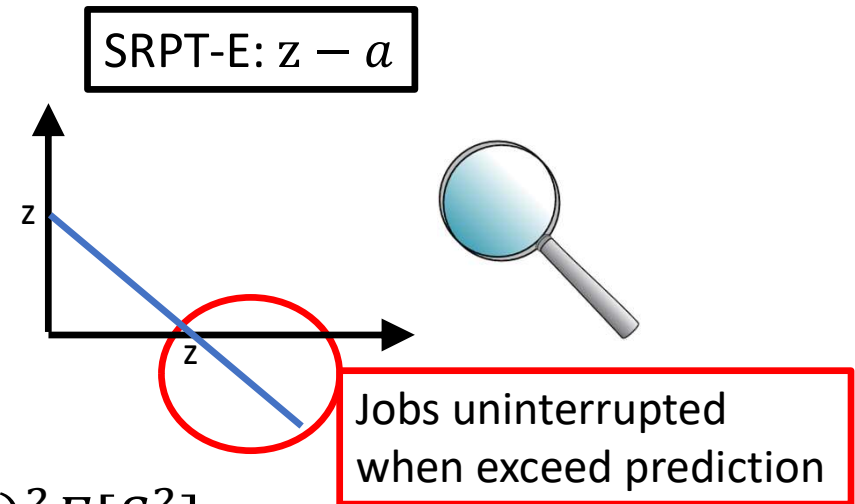
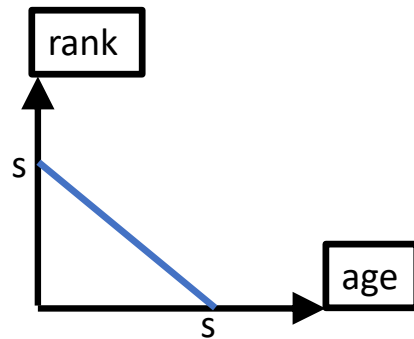
Consistency: $\lim_{\alpha, \beta \rightarrow 1} \frac{E[T^\pi]}{E[TSRPT]} = 1$

Graceful Degradation: $\forall \alpha, \beta, \frac{E[T^\pi]}{E[TSRPT]} \leq c \frac{\alpha}{\beta}$

Robustness: $\forall \alpha, \beta, \frac{E[T^\pi]}{E[TSRPT]} \leq d$

Naïve Scheduling Policy

SRPT:
Rank = $s - a$
Lower is better

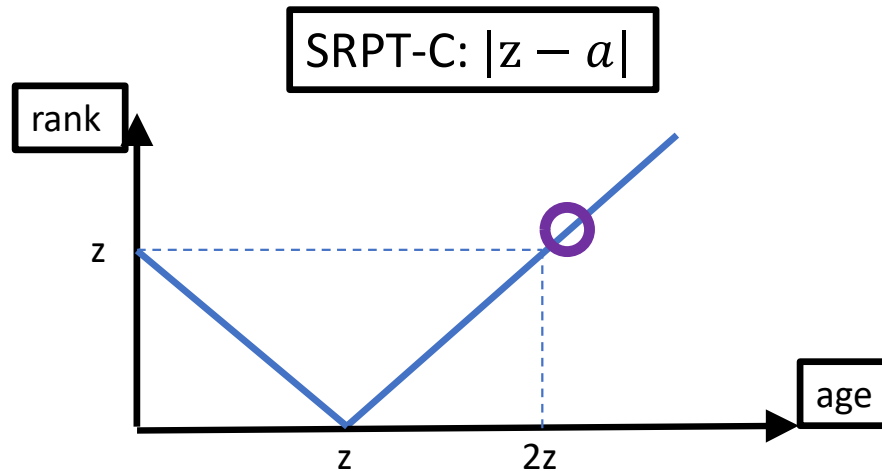


Lower bound: $E[T^{SRPT-}] \geq \frac{\lambda}{2} (1 - \beta)^2 E[S^2]$

If $E[S^2] = \infty, E[T] = \infty$.

 In datacenter and webserver traces, $E[S^2]$ very large

SRPT with Checkmark

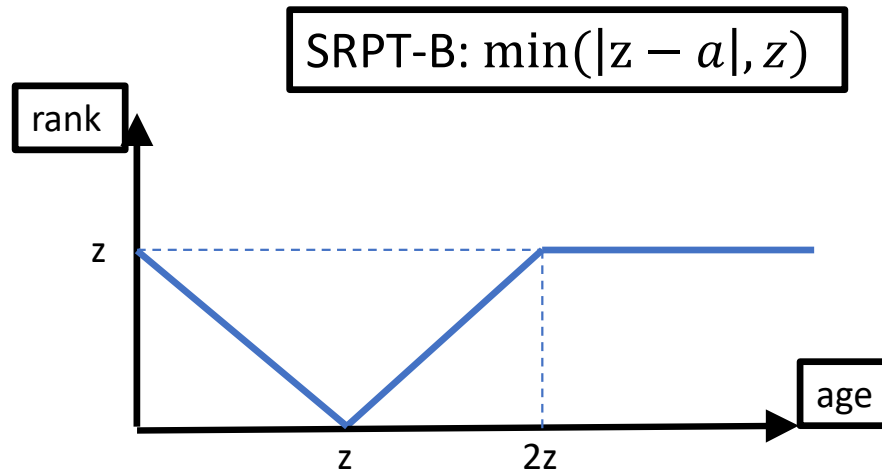


SRPT-C is consistent: $\lim_{\alpha, \beta \rightarrow 1} \frac{E[T^{SRPT-}]}{E[T^{SRPT}]} = 1$

First consistent policy!

Not gracefully degrading: if $\beta < \frac{1}{2}$, $\frac{E[T^{SRPT-}]}{E[T^{SRPT}]}$ unbounded.

SRPT with Bounce



First policy to be consistent and gracefully degrading!



Results

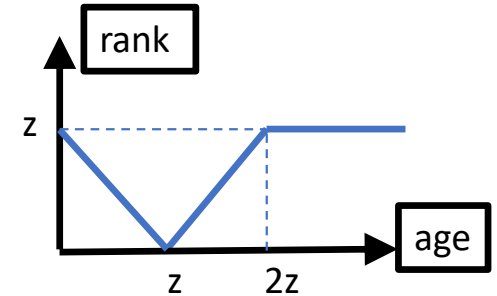
$$\frac{E[T^{SRPT}]}{E[T]} \leq \frac{\alpha}{\beta} K(\alpha, \beta)$$

where $K(\alpha, \beta) = 1 + \left(\frac{3}{2}\alpha 1\{\beta < 1\} + 1\right) \min\left\{1, \max\left\{1 - \frac{1}{\alpha}, \frac{1}{\beta} - 1\right\}\right\}$

$\lim_{\alpha, \beta \rightarrow 1} K(\alpha, \beta) = 1$ (Consistency)

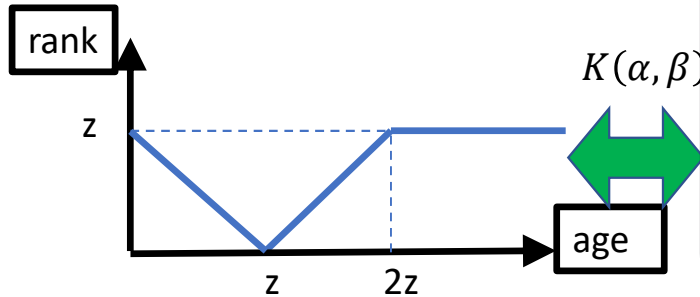
$\forall \alpha, \beta, K(\alpha, \beta) \leq 3.5$ (Graceful degradation)

Robustness impossible, Gittins policy lower bound

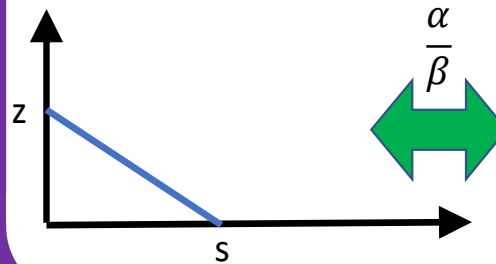


Proof Methods

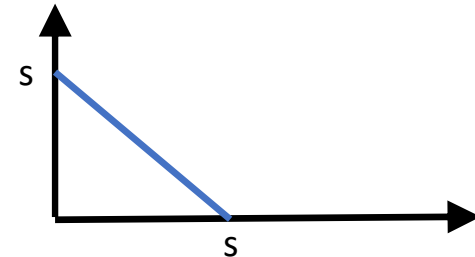
$$\text{SRPT-B: } \min(|z - a|, z)$$



$$\text{SRPT-SE: } \frac{z}{s} (s - a)$$



$$\text{SRPT: } s - a$$



Proof uses cutting-edge queueing theory:

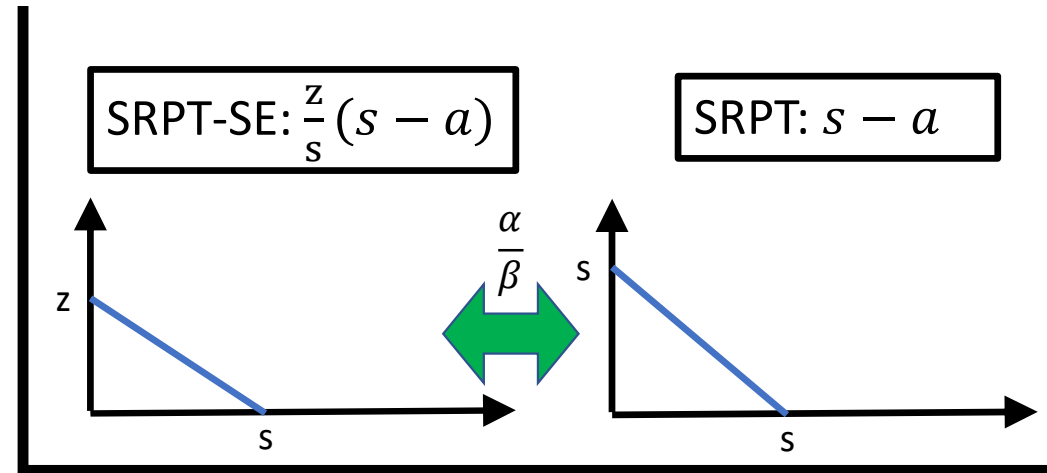
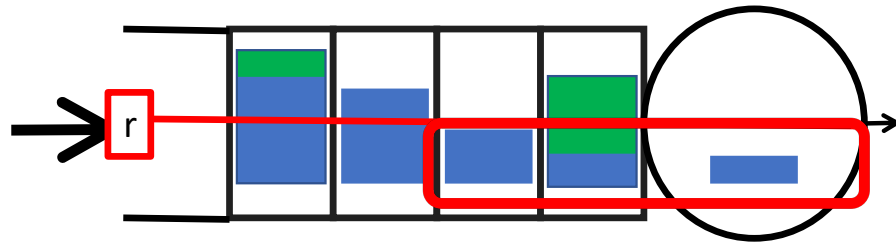
SOAP

Analysis of $E[T]$ for rank-function policies

WINE

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

SRPT-SE vs. SRPT



r-relevant work: $W(r) = \sum (s_i - a) 1\{s_i - a \leq r\}$

Thm: SRPT minimizes $W(r)$: $E[W^{SRPT}(r)] \leq E[W^\pi(r)] \forall \pi, r$

Thm: SRPT-SE nearly min. $W(r)$: $E[W^{SRPT-SE}(r)] \leq E\left[W^\pi\left(\frac{\alpha}{\beta}r\right)\right] \forall \pi, r$

$$E[W^{SRPT-SE}(r)] \leq E\left[W^{SRPT}\left(\frac{\alpha}{\beta}r\right)\right]$$

$$E[T^{SRPT-S}] \leq \frac{\alpha}{\beta} E[T^{SRPT}]$$

$$\text{WINE: } E[T] = \frac{1}{\lambda} \int_0^\infty \frac{E[W(r)]}{r^2} dr$$

Zooming Out on Scheduling with Predictions

Today: Unknown prediction distribution (S, Z)

Known prediction distribution:

- Single server solved by Gittins policy

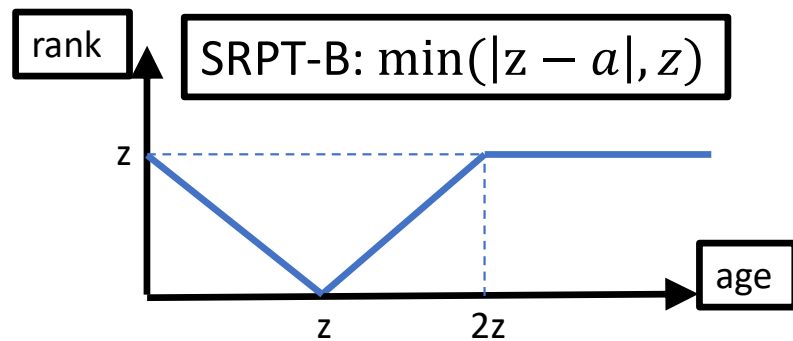
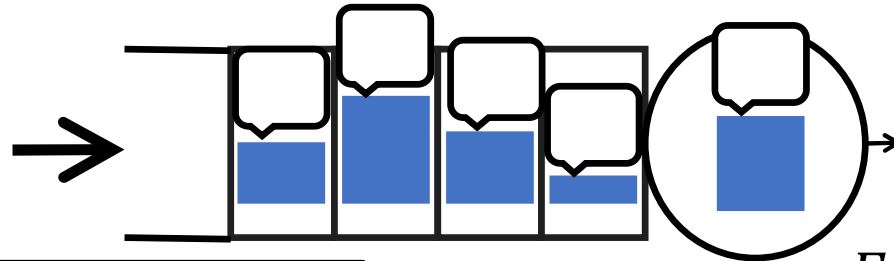
- Multiserver: “The Gittins Policy is Nearly Optimal in the M/G/k under Extremely General Conditions”, [G., Scully, Harchol-Balter], SIGMETRICS 2021

Strategic predictions:

- “Incentive Compatible Queues Without Money”, [G., Mitzenmacher], arXiv

Stochastic Scheduling with Predictions

igros@cmu.edu
isaacg1.github.io



$$\lim_{\alpha, \beta \rightarrow 1} \frac{E[T^{SRPT}]}{E[T^{SRPT}]} = 1$$

$$\frac{E[T^{SRPT-}]}{E[T^{SRPT}]} \leq 3.5 \frac{\alpha}{\beta}$$

