

Lottery Scheduling

I. Lottery Scheduling

Very general, proportional-share scheduling algorithm.

Problems with traditional schedulers:

- o Priority systems are ad hoc at best: highest priority always wins
- o “Fair share” implemented by adjusting priorities with a feedback loop to achieve fairness over the (very) long term (highest priority still wins all the time, but now the Unix priorities are always changing)
- o Priority inversion: high-priority jobs can be blocked behind low-priority jobs
- o Schedulers are complex and difficult to control

Lottery scheduling:

- o Priority determined by the number of tickets each process has: priority is the relative percentage of all of the tickets competing for this resource.
- o Scheduler picks winning ticket randomly, gives owner the resource
- o Tickets can be used for a wide variety of different resources (uniform) and are machine independent (abstract)

How fair is lottery scheduling?

- o If client has probability p of winning, then the expected number of wins (from the binomial distribution) is np .
- o Variance of binomial distribution: $\sigma^2 = np(1-p)$
- o Accuracy improves with \sqrt{n}
- o Geometric distribution used to find tries until first win
- o Big picture answer: mostly accurate, but short-term inaccuracies are possible; see Stride scheduling below.

Ticket Transfer: how to deal with dependencies

- o Basic idea: if you are blocked on someone else, give them your tickets
- o Example: client-server
 - Server has no tickets of its own
 - Clients give server all of their tickets during RPC
 - Server’s priority is the sum of the priorities of all of its active clients
 - Server can use lottery scheduling to give preferential service to high-priority clients
- o Very elegant solution to long-standing problem (not the first solution however)

Scheduling

Ticket inflation: make up your own tickets (print your own money)

- o Only works among mutually trusting clients
- o Presumably works best if inflation is temporary
- o Allows clients to adjust their priority dynamically with zero communication

Currencies: set up an exchange rate with the base currency

- o Enables inflation just within a group
- o Simplifies mini-lotteries, such as for a mutex

Compensation tickets: what happens if a thread is I/O bound and regular blocks before its quantum expires? Without adjustment, this implies that thread gets less than its share of the processor.

- o Basic idea: if you complete fraction f of the quantum, your tickets are inflated by $1/f$ until the next time you win.
- o Example: if B on average uses $1/5$ of a quantum, its tickets will be inflated 5x and it will win 5 times as often and get its correct share overall.
- o What if B alternates between $1/5$ and whole quanta?

Problems:

- o Not as fair as we'd like: mutex comes out 1.8:1 instead of 2:1, while multimedia apps come out 1.92:1.50:1 instead of 3:2:1
- o Practice midterm question: are these differences statistically significant? (probably are, which would imply that the lottery is biased or that there is a secondary force affecting the relative priority)
- o Multimedia app: biased due to X server assuming uniform priority instead of using tickets. Conclusion: to really work, tickets must be used everywhere. Every queue is an implicit scheduling decision... Every spinlock ignores priority...
- o Can we force it to be unfair? Is there a way to use compensation tickets to get more time, e.g., quit early to get compensation tickets and then run for the full time next time?
- o What about kernel cycles? If a process uses a lot of cycles indirectly, such as through the ethernet driver, does it get higher priority implicitly? (probably)

Stride Scheduling: follow on to lottery scheduling (not in paper)

- o Basic idea: make a deterministic version to reduce short-term variability
- o Mark time virtually using "passes" as the unit
- o A process has a stride, which is the number of passes between executions. Strides are inversely proportional to the number of tickets, so high priority jobs have low strides and thus run often.
- o Very regular: a job with priority p will run every $1/p$ passes.
- o Algorithm (roughly): always pick the job with the lowest pass number. Updates its pass number by adding its stride.
- o Similar mechanism to compensation tickets: if a job uses only fraction f , update its pass number by $f \times stride$ instead of just using the stride.
- o Overall result: it is far more accurate than lottery scheduling and error can be bounded

Scheduling

absolutely instead of probabilistically