

گزارش امتیازی (ALFS (Heap-based CFS Simulation))

این گزارش سه محور دارد:

1. چرا CFS واقعی از RB-Tree استفاده می‌کند و Heap چه trade-off هایی دارد.
2. جایگزین‌های مطرح CFS (بهویژه EEVDF) و مقایسه‌ی پژوهشی.
3. اثر معماری‌های asymmetric / big.LITTLE روی زمان‌بندها و implication طراحی.

یادآوری: ALFS در این پژوهه یک شبیه‌ساز user-space است و دقیقاً کرنل لینوکس نیست.

1) RB-Tree در برابر CFS در «همه‌چیز فقط» extract-min نیست

1.1 CFS دقیقاً از RB-Tree استفاده می‌کند؟

مستند طراحی CFS در کرنل لینوکس صریحاً می‌گوید runqueue یک rbtree مرتب شده بر اساس vruntime نگه می‌دارد و «leftmost» (کمترین vruntime) را برای اجرا انتخاب می‌کند. ایده‌ی کلیدی هم این است که با پیشرفت زمان، entity‌های اجراسده آرام‌آرام به سمت راست می‌روند تا هر task در یک بازه‌ی قابل‌پیش‌بینی، فرصت CPU بگیرد. منبع: [K1].

1.2 آیا Heap از نظر انتخاب min بهتر نیست؟

Heap peek-min فقط به CFS عالی است؛ اما در کرنل واقعی محدود نیست. چند peek-min برای عملیات مهم دیگر:

A) task حذف/بروزرسانی یک (arbitrary removal/update)

- در کرنل، task‌ها مرتب unblock/block می‌شوند، از runqueue حذف می‌شوند، و دوباره insert می‌شوند.
- در RB-tree: با داشتن اشاره‌گر به node، حذف $O(\log n)$ است.
- در heap ساده اگر index نگه ندارید، حذف می‌تواند $O(n)$ شود چون باید عنصر را پیدا کنید.
- (در ALFS ما برای حل این مشکل heap داخل index را نگه می‌داریم و حذف را $O(\log n)$ می‌کنیم.)
- این نکته در بحث‌های فنی هم مطرح شده که heap بدون "handle" برای حذف دلخواه در دست دارد. [K2]

B) traversal مرتب (in-order iteration)

در برخی مسیرها/heuristic، داشتن ability برای حرکت "به ترتیب" در range (successor/predecessor) یا کار با مفید است.

- RB-tree به شکل طبیعی این قابلیت را می‌دهد.

- کامل را نمی‌دهد order "کمترین" را به شکل مستقیم می‌دهد و

C) SMP / load-balancing

در لینوکس واقعی، زمان‌بندی per-CPU است و load balancing یک بخش مهم طراحی است. داشتن یک ساختار مرتب می‌تواند طراحی و پیاده‌سازی heuristic‌های انتخاب/مهاجرت را ساده‌تر کند. [K1]

D) group scheduling (cgroup)

گروه‌های بیشتری داریم و عملیات‌های جانبی بیشتر می‌شود. RB-tree به خاطر group scheduling، entity traversal و operations می‌تواند استاندارد، انتخاب قابل دفاعی است.

1.3 Heap است؟ کجا بهتر

- است و پیاده‌سازی اش می‌تواند ساده و سریع باشد cache-friendly آرایه‌ای heap.
- `extract-min` و `insert` در $O(\log n)$ و `peek-min` در $O(1)$.

اما هزینه‌ی workload در heap کمتر است بالا می‌رود:

- چون heap به شما "کمترین eligible برای k CPU" را نمی‌دهد، مجبور می‌شوید pop/push موقت انجام دهید (بدترین حالت نزدیک به $O(n)$ per tick در سناریوهای pathological مثل affinity mask های محدود).

نکته‌ی عملی: در این پروژه workload‌ها کوچک‌اند و این worst-case عملیات را کمتر خ می‌دهد؛ اما در کرنل واقعی با هزاران task و topology پیچیده، این هزینه می‌تواند مهم شود.

1.4 نتیجه‌گیری بخش ۱

بلکه برای مجموعه‌ای از `min` یک انتخاب "سیستمی" است: نه فقط برای انتخاب CFS برای RB-tree، بلکه برای انتخاب "سیستمی" است: نه فقط برای انتخاب CFS برای RB-tree، قابل استفاده است، اما هزینه‌ی overhead eligible، insert/traversal،Heap،SMP،گروه‌ها،زیاد این هزینه می‌تواند به سازد.

2) جایگزین‌های CFS و تمرکز روی EEVDF

2.1 EEVDF چیست و چرا مهم است؟

است که `EEVDF` (Earliest Eligible Virtual Deadline First) یک الگوریتم proportional share ایده‌ی اصلی اش این است:

- هر entity یک مفهوم virtual time دارد،
- ولی انتخاب بر اساس "earliest eligible virtual deadline" انجام می‌شود، نه صرفاً کمترین `.vruntime`

دو مرجع پژوهشی مهم:

- Ion Stoica, Hussein Abdel-Wahab, "Earliest Eligible Virtual Deadline First: A Flexible and Accurate Mechanism for Proportional Share Resource Allocation", TR-95-22, 1995. [E1]
- Scott A. Brandt, James B. Weissman, "EEVDF Proportional Share Resource Allocation Revisited", RTSS 2000 (PDF proceedings). [E2]

2.2 تفاوت ایده‌ای CFS با EEVDF

- را اجرا می‌کند "کار می‌کند و "کمترین vruntime عمدتاً با CFS.
- EEVDF دارد virtual time، eligible time و virtual deadline مفهوم علاوه بر:
 - entity ها تا شوند وارد رقابت اصلی نیستند" earliest virtual deadline با entity سپس می‌شود.

در لینوکس هم یک مستند رسمی درباره‌ی EEVDF وجود دارد که همین مفاهیم را برای context کرنل توضیح می‌دهد. [K3]

نکته‌ی مهم: EEVDF صرفاً یک ایده‌ی پژوهشی نیست؛ لینوکس در سال‌های اخیر آن را در چارچوب CFS وارد کرده است تا latency و fairness بهتر و دقیق‌تر به دست آید. [K3]

2.3 جایگزین‌های دیگر (BFS / MuQSS و ...)

برای یک منبع. BFS و MuQSS (Con Kolivas interactivity/desktop) مطرح شدند thesis که استفاده [M1] را مقایسه کرده، می‌توانید از MuQSS و CFS در حد) دانشگاهی قابل استناد کنید.

2.4 ALFS جایگاه

در این پژوهه:

- از "فلسفه‌ی CFS و وزن nice" را نگه می‌دارد،
- اما ساختار داده‌ی انتخاب را از RB-tree به heap تغییر می‌دهد،
- و (با cgroup) یک مدل ساده از hierarchical fairness را شبیه‌سازی می‌کند.

پس ALFS بیشتر یک آزمایش مهندسی روی ساختار داده است تا یک الگوریتم scheduling جدید مثل .EEVDF

3) big.LITTLE CPU روی زمان‌بندها fairness: چه؟

3.1 مشکل اصلی big.LITTLE

:big.LITTLE در

- coreها capacity. یکسان ندارند

- یک tick روی big با یک tick روی LITTLE برابر نیست (از نظر "کار انجام شده").
- migration latency را خراب کند و میتواند cache/warmup بین خوشها هزینه‌ی

دو مرجع پژوهشی مفید:

- Vinicius Petrucci, Orlando Loques, "Lucky scheduling for energy-efficient heterogeneous multi-core systems", USENIX HotPower 2012. [B1]
- Agostino Mascitti et al., "Dynamic partitioned scheduling of real-time tasks on ARM big.LITTLE architectures", Journal of Systems Architecture (2020). [B2]

3.2 پیام برای طراحی زمان‌بند

A) fairness باید capacity-aware شود

اگر CPU‌ها همسان نباشند، fairness صرفاً بر حسب زمان (time fairness) می‌تواند غلط‌انداز باشد. بعضی طرح‌ها "work-based" یا "capacity-normalized" را پیشنهاد می‌کنند.

B) energy-aware placement

در لینوکس، Energy Aware Scheduling (EAS) برای topologies ناهمسان اضافه شده و از مدل انرژی برای تصمیم placement استفاده می‌کند. [K4][K5]

C) overhead مهاجرت

بسیاری از سیاست‌ها تلاش می‌کنند migration را کاهش دهند تا latency و thrashing را کم شود.

3.3 پیشنهادهای طراحی ALFS برای big.LITTLE-aware بخواهیم شویم (پیشنهادی)

(در این پژوهه پیاده نشده، ولی برای دفاع/جمع‌بندی پژوهشی مفید است)

1. تعريف `:capacity` CPU بر اساس `delta_exec`

- `effective_exec = quanta_us * cpu_capacity[cpu]` سپس `effective_exec` update بر اساس vruntime شود.

2. penalty: برای مهاجرت

• اگر task CPU قبلی eligible است، ترجیح به sticky CPU بده.

• 3. برای quota/shares گروهها، accounting را بر اساس "work" time "ne صرفاً" انجام دهید.

نتیجه‌گیری: آیا ALFS «بهتر» است؟

- برای شبیه‌سازی ساده و آموزش ساختار داده، heap قابل قبول و حتی جذاب است.
- برای کرنل واقعی با SMP، load balancing و عملیات‌های جانبی زیاد، RB-tree انتخابی قابل دفاع است. [K1]
- الگوریتم‌هایی مثل EEVDF تغییر مفهومی عمیق‌تری نسبت به "تعویض ساختار داده" هستند و [K3][E2][E1] بهتر با حفظ latency هدفشان proportional fairness است.

References

Kernel / Docs

- [K1] Linux Kernel Documentation — CFS design: [sched-design-CFS](#)
<https://www.kernel.org/doc/html/latest/scheduler/sched-design-CFS.html>
- [K2] دربارهی RB-tree vs heap (برای CFS) بحث فنی جمع‌بندی غیررسمی :
نہ منبع (رسمی کرنل):
<https://stackoverflow.com/questions/33191110/reason-why-cfs-scheduler-using-red-black-tree>
- [K3] Linux Kernel Documentation — EEVDF scheduler:
<https://www.kernel.org/doc/html/latest/scheduler/sched-eevdf.html>
- [K4] Linux Kernel Documentation — Energy Aware Scheduling (EAS):
<https://www.kernel.org/doc/html/latest/scheduler/sched-energy.html>
- [K5] LWN — Energy Aware Scheduling overview (زمنیه/تاریخچه):
<https://lwn.net/Articles/749738/>

EEVDF (papers)

- [E1] Ion Stoica, Hussein Abdel-Wahab (1995), TR-95-22 (PDF):
<https://people.eecs.berkeley.edu/~istoica/papers/eevdf-tr-95.pdf>
- [E2] Scott A. Brandt, James B. Weissman (2000), RTSS proceedings paper (PDF):
<https://users.soe.ucsc.edu/~sbrandt/rtss2000/proceedings/6.pdf>

big.LITTLE / Heterogeneous scheduling (papers)

- [B1] Petrucci, Loques (2012), HotPower (PDF):
<https://www.usenix.org/system/files/conference/hotpower12/hotpower12-final34.pdf>
- [B2] Mascitti et al. (2020), Journal of Systems Architecture (ScienceDirect page):
<https://www.sciencedirect.com/science/article/pii/S0164121220302764>

Other scheduler comparisons

- [M1] “CFS and MuQSS Comparison” (DiVA portal thesis PDF):
<https://www.diva-portal.org/smash/get/diva2:1566002/FULLTEXT01.pdf>