

# ALFS — Anushiravan-Level Fair Scheduler (User-space)

این مخزن یک  **شبیه‌ساز User-space** از زمان‌بند CFS است: نسخه‌ای تغییریافته از ایده‌ی استفاده می‌کند و از طریق  **Unix Domain Socket** روابط را به جای RB-Tree به صورت  **JSON** دریافت کرده و برای هر  **vtime** یک خروجی زمان‌بندی  **JSON** برمی‌گرداند.

نکته: این پروژه «کرنل» نیست؛ ولی سعی کرده مفاهیم کلیدی CFS (vruntime + وزن + nice + fairness) را با یک مدل پیاده کند.

## Build

make

خروجی: **alfs/.**

## Run

پارامترها:

- **--cpus N** تعداد CPU ها
  - واحد: میکروثانیه (microsecond) tick کوانتای هر **quanta-us Q**
  - (اختیاری) **socket PATH** مسیر UDS
  - اگر ندهید، برنامه اول **event.socket/.** و اگر نبود **socket.event/.** را امتحان می‌کند.
  - اگر هیچ‌کدام وجود نداشت/سرور بالا نبود، با پیام **connect: No such file or directory** و کد خروج **1** خارج می‌شود.
- **--burst-mode MODE** (اختیاری)
  - **freeze\_pushback** هنگام (پیش‌فرض) **CPU\_BURST** ابتدا **vruntime = max\_vruntime** سپس فریز می‌شود.
  - **freeze**: فقط **vruntime** (بدون pushback)
  - **track**: **vruntime** کلاً **update** مثل حالت عادی **hint** نادیده گرفتن) می‌شود
- **--meta-minimal** فقط **{preemptions,migrations}** meta را در خروجی می‌دهد (پیش‌فرض)
- **--meta-extra** علاوه بر **{preemptions,migrations}** این‌ها را هم اضافه می‌کند:
  - **runnableTasks**
  - **blockedTasks**
- **--debug** لگ‌های داخلی (stderr)

نمونه:

```
./alfs --cpus 4 --quanta-us 1000 --socket socket.event --debug
```

---

## I/O Protocol (UDS + JSON framing)

- برنامه به عنوان **CLIENT** به UDS وصل می‌شود.
- ورودی: برای هر **vtime** یک JSON دریافت می‌شود (TimeFrame).
- خروجی: برای همان **vtime** یک JSON در قالب SchedulerTick تولید می‌شود و با newline خاتمه پیدا می‌کند.

## Robust framing

چون brace/bracket دقیق تضمین نشده، ورودی به صورت stream خوانده می‌شود و با یک JSON کامل استخراج می‌کنیم:

- شروع از { یا ]
- شمارش عمق { و [ (با در نظر گرفتن escape و string)
- وقتی depth به ۰ رسید، یک JSON کامل داریم.

این روش هم با newline-delimited newline و هم بدون newline جواب می‌دهد.

---

## Data Model

### Task

هر ترسک:

- **id** (string)
- **nice** [19+..20-] در بازه‌ی
- **weight** از جدول **prio\_to\_weight** hardcode شده (داخل کد) کرنل
- **vruntime\_us** (virtual runtime «virtual microseconds»)
- **state: RUNNABLE | BLOCKED | EXITED**
- **cpuMask** (bitset)
- **cgroup** (pointer)
- **burst\_remaining** و **burst\_freeze** برای CPU\_BURST

### Cgroup

:cgroup هر

- **cpuShares** وزن گروه
- **cpuQuotaUs** و **cpuPeriodUs** برای throttling
- **cpuMask** (bitset)

- `vruntime_us` گروه
  - `task_heap`: min-heap داخل گروه runnable تسک‌های (کلید) `task->vruntime_us` خود هم داخل cgroup می‌گیرد (کلید) `cg->vruntime_us` (سپس از heap در میان گروه‌ها pop می‌کنیم) •
- 

## Scheduling Algorithm (ELIAD: Explain Like I'm Defending It)

### ایده‌ی CFS که نگه داشتیم

- CFS می‌خواهد هر تسک متناسب با وزنش CPU بگیرد.
- معیار ساده‌ی پیشرفت: `vruntime`
- هر بار که تسک روی CPU اجرا می‌شود:
- `vruntime += delta_exec * NICE_0_WEIGHT / weight`
- پس تسکی که کمتر CPU گرفته، `vruntime` کمتر دارد و زودتر انتخاب می‌شود.
- در این پروژه:
- `delta_exec = quanta_us` کوانتای هر tick)

### چرا Heap؟

در CFS واقعی از RB-Tree استفاده می‌شود؛ ما اینجا به جای tree، heap را انجام می‌دهیم.

## چطور اضافه شد؟ Group Scheduling (cgroup)

دو سطح :heap

1. `cg_heap` (heap گروه با کمترین گروهها: `cg->vruntime_us`) را انتخاب می‌کند
2. داخل همان گروه، از `task->vruntime_us` کمترین `task_heap` انتخاب می‌شود

update ها:

- وقتی یک task از گروه روی CPU اجرا شد:
- `cg->vruntime_us += quanta_us * NICE_0_WEIGHT / cg->cpuShares`
- `task->vruntime_us += quanta_us * NICE_0_WEIGHT / task->weight` (مگر CPU\_BURST)

### CPU + Affinity چند

برای هر CPU به ترتیب `cpu=0..N-1`

- از `cg_heap` pop می‌کنیم تا گروهی پیدا شود که:
- throttled نباشد
- runnable task داشته باشد
- cpuMask گروه آن CPU را اجازه بدهد
- سپس از heap داخلی گروه pop می‌کنیم تا taskی پیدا شود که:

- cpuMask بدهد CPU خودش را اجازه دهد
  - در همین tick روی CPU دیگر انتخاب نشده باشد (`current_cpu == -1`)
  - task برزنمی‌گردد تا دوبار انتخاب نشود heap دوباره به tick انتخاب شده تا پایان task push گروهشان می‌شوند.
  - بعد از اینکه همه CPU‌ها انتخاب شدند، task‌ها دوباره به heap گروهشان push می‌شوند.
  - این کار deterministic است چون:
  - انتخاب CPU‌ها ترتیب ثابت دارد
  - tie-break در heap با `(vruntime, seq, id)` است
- 

## Event Semantics

- `TASK_CREATE`: ایجاد + ورود به گروه runnable heap
- `TASK_BLOCK`: و رفتن به heap و خروج از BLOCKED
- `TASK_UNBLOCK`: با همان heap بازگشت به vruntime
- `TASK_YIELD`: `vruntime = max_vruntime` در heap
- `TASK_SETNICE`: فعلی دست نمی‌خورد vruntime وزن بعض می‌شود!
- `TASK_SET_AFFINITY`: cpuMask عوض می‌شود
- `CGROUP_*`: (رد می‌شود non-empty گروه) ساخت/تغییر/حذف گروه
- `TASK_MOVE_CGROUP`: جدید اضافه می‌شود runnable heap حذف و به باشد از duration

## CPU\_BURST (تفسیر پیاده‌سازی)

- به صورت پیش‌فرض:
- وقتی `CPU_BURST(taskId, duration)` آمد:
  - `task->vruntime` گروه می‌بریم (pushback) اگر روی CPU اجرا شد، `duration` tick سپس تا `vruntime update` نمی‌شود
  - بعد از تمام شدن duration، دوباره `vruntime update` می‌شود.
  - می‌توانید این را با `burst-mode freeze|track--` تغییر دهید.
- 

## Metadata دقیقاً چیست؟

- در خروجی همیشه این دو مقدار وجود دارد:
- `idle`/دیگری task `tick` را قبل رویشان یک `tick` که task `preemptions` تعداد دارد: قرار گرفته CPU روی آن.
  - نکته: تغییر `idle -> task` را `idle` حساب نمی‌کنیم (همان الگوی مثال پروژه).
  - قبلی‌شان این CPU قبل تغییر کرده و `task` شان نسبت به CPU که task `migrations` تعداد شده است.

- این تعریف طوری انتخاب شده که swap را بین دو task حساب نکند (مثل مثال پروژه).

پیشفرض، خروجی **minimal** است. اگر **meta-extra--** بدهید، این دو لیست هم اضافه می‌شود:

- **runnableTasks**
  - **blockedTasks**
- 

## Debugging

```
./alfs --cpus 4 --quanta-us 1000 --debug
```

چیزهایی که چاپ می‌شود:

- ساخت task/cgroup
  - انتخاب task برای هر CPU در هر tick
- 

## Tests (mini harness)

یک سرور ساده در **tests/server.py** وجود دارد که:

- یک UDS server ایجاد می‌کند
- فایل trace (مثل **tests/trace\_demo.jsonl**) را به عنوان line-by-line می‌خواند
- هر خط را به عنوان یک TimeFrame می‌فرستد
- پاسخ ALFS را می‌خواند و چاپ می‌کند

### اجرای دقیق (دو ترمیناله)

ترمینال ۱ (سرور):

```
rm -f socket.event event.socket
python3 tests/server.py socket.event
tests/trace_demo.jsonl`
```

ترمینال ۲ (کلاینت ALFS):

```
./alfs --cpus 1 --quanta-us 1000 --meta-minimal --socket socket.event
```

برای تست‌های چند-CPU و meta کامل:

```
./alfs --cpus 2 --quanta-us 1000 --meta-extra --socket socket.event
```

---

## Files

- **src/alfs.c**: همه چیز (heap + scheduler + JSON + UDS)
  - **third\_party/jsmn.h**: JSON parser کوچک
  - **tests/server.py**: UDS سرور تست
  - **tests/trace\_demo.jsonl**: یک سناریو ساده
  - **report/REPORT.md**: گزارش امتیازی (RB-Tree vs Heap + EEVDF + big.LITTLE)
- 

## ”راهنمای دفاع“ (Bullet های طلایی)

۱) مدل ذهنی: «CFS» تلاش می‌کند CPU را متناسب با وزن تقسیم کند؛ vruntime نشان می‌دهد چه کسی کمتر سهم گرفته.»

۲) معادله **:vruntime**

**vruntime += delta\_exec \* NICE\_0\_WEIGHT / weight**

(در این شبیه‌ساز: **delta\_exec = quanta\_us**)

۳) چرا **heap** کار می‌کند: چون انتخاب بعدی تقریباً همیشه "کمترین" vruntime است  $\Rightarrow$  **extract-min**

۴) مشکل **heap** نسبت به **RB-tree** (برای گزارش): حذف دلخواه / traversal مرتب / عملیات load و "یافتن نزدیکترین‌ها" balancing / group scheduling راحت‌تر است.

۵) چند **CPU + affinity**: چون heap min فقط CPU pop می‌کنیم تا به گزینه‌ی eligible برسیم، موارد نامناسب را موقت کنار می‌گذاریم و بر می‌گردانیم.

۶) **cgroup**: heap (group task)، vruntime و سپس shares scale گروه با سطح دو شود.

۷) **quota/period**: runtime period در هر گروه رد شد quota جمع می‌شود، وقتی از period throttled و تا اجازه‌ی اجرا ندارد period پایان.

۸) **CPU\_BURST**: policy شفاف + قابل تغییر با flag.

۹) **determinism**: tie-break مشخص و انتخاب CPU ثابت ترتیب ها به می‌شوند.