# Sound of Scheduling

## Writing Linux Schedulers in Java

Johannes Bechberger

David Kiefer

Sport

Cook

Sleep

Work

Task 1 Task 2

CPU

1

Time

Hear this sound?

It's my scheduler

It's my scheduler

Written in Java

# Why?

"The only way of discovering the limits of the possible is to venture a little way past them into the impossible.

Clarke's second law

# How to modify the kernel?

# Traditional ways

1. Change the Kernel
2. Kernel module

# Traditional ways

1. Change the Kernel
2. Kernel module

Interfaces are complicated

# Traditional ways

1. Change the Kernel
2. Kernel module

Not possible with schedulers

# Problem: Only a few are implemented on your system

# CFS was built in a simpler time

- Much smaller CPUs
- Topologies much more homogeneous
- Cores spaced further apart, migration cost typically high
- Power consumption and die area wasn't as important
- The fundamental assumptions behind heuristics may be easier to justify

Just two cores

Just one L3 cache



Intel Xeon MP 71xx die

# Architectures *much* more complicated now

- Heterogeneity is becoming the norm
- Non-uniform memory accesses between sockets
- Non-uniform memory accesses between CCDs
- Non-uniform memory accesses between CCXs
- Non-uniform memory accesses between CCXs in the same CCD



AMD Zen 2 Rome

4 cores per "CCX"
8 cores per "CCD"

2 L3 caches per CCD!

8 cores per "CCX"
8 cores per "CCD"

1 L3 cache per CCD!

# Let's create our own

# Let's create our own

Has someone done this before in this room?

# Let's create our own

*How hard can it be?*

# Let's create our own

# Who was there in the Firewall Talk?

# Skip ahead

"

eBPF is a crazy technology, it's like putting JavaScript into the Linux kernel

Brendan Gregg

> " eBPF is a crazy technology, it's like putting JavaScript into the Linux kernel

Brendan Gregg

# eBPF runtime

# eBPF runtime



source code → eBPF program

Compilation

bytecode → eBPF program

Development

bpf Syscall

Verifier

Linux Kernel

Courtesy of Mohammed Aboullaite

Courtesy of Mohammed Aboullaite

# How to share data?



via sockets:

Program A → sendMessage( · ) → Program B

Data(uid=0,
     gid=0,
     counter=10)

via shared memory:

Program A —set→ [Data(uid=0, gid=0, counter=10)] ←get— Program B

# How to share data?



via sockets:

Data(uid=0,
     gid=0,
     counter=10)

Program A

sendMessage(    ·    )

Program B

## Any Problems?

via shared memory:

Program A

set

Data(uid=0,
     gid=0,
     counter=10)

get

Program B

# How to share data?

via eBPF maps:

# eBPF Maps

- I want to use a programming language
which doesn't only run in Windows.

- I want to use a programming language
which doesn't only run in user land

# Demo

# XDP

# Back to scheduling

sched_ext

# Sched Ext

The extensible sched_class

P. F. C. L.

Penguins For Cache Locality

Will work for CPU cycles

David Vernet
Kernel engineer

∞ Meta

> 1. Ease of experimentation and exploration

> 2. Customization

> 3. Rapid scheduler deployments

# Typical Scheduler Goals

# Fairness

Typical Scheduler Goals

# Resource Utilization

# Typical Scheduler Goals

# Overhead

# Typical Scheduler Goals

# Responsiveness

Example Application

sched_ext

eBPF    hello eBPF

# Let's create a scheduler

```java
@BPF(license = "GPL")
abstract class SampleScheduler
    extends BPFProgram
    implements Scheduler, Runnable {
    // ...
}
```

```
PID             Process Name              Enqueue Count
----------------------------------------------------------
204358          java                                  102
204403          ForkJoinPool.co                        78
204406          ForkJoinPool.co                        76
204407          ForkJoinPool.co                        75
204402          ForkJoinPool.co                        74
204399          ForkJoinPool.co                        72
204404          ForkJoinPool.co                        71
204412          ForkJoinPool.co                        70
204405          ForkJoinPool.co                        69
204401          ForkJoinPool.co                        68
```

# What is the performance?

# Good*

*\* For a typical Java benchmark*

# How does it work?

# Let's see some schedulers

# First-Come, First-Served Scheduler

*Run as long as you want, we won't stop you*

# FCFSScheduler

# Making errors is normal

```
/**
 * @timeout_ms: The maximum amount of time, in milliseconds, that a
 * runnable task should be able to wait before being scheduled. The
 * maximum timeout may not exceed the default timeout of 30 seconds.
 *
 * Defaults to the maximum allowed timeout value of 30 seconds.
 */
u32 timeout_ms;
```

https://github.com/torvalds/linux/blob/master/kernel/sched/ext.c

# First-Come, First-Out Scheduler

*The early bird eats the time slice*

# MinimalScheduler

# Lottery Scheduler

*Are you the lucky task who gets the time slice?*

# LotteryScheduler

# VRuntime-based Scheduler

- Tracks virtual runtime (vruntime) of tasks (time on CPU)
- Task with shortest vruntime runs first
- Use a simple priority queue

runtime

# VRuntime-based Scheduler

*You already run quite a long time, lets choose another task*

# Proportional weight-based CPU allocation: fairness

- Each task $T_i$ has a weight $w_i$
- The runtime assigned to each task $T_i$ is proportional to its weight $w_i$ divided by the sum of all the runnable tasks' weight

$$runtime(T_i) = \int_{t_0}^{t_1} \frac{w_i}{\sum_{j=0}^{N} w_j} dt \simeq \frac{w_i}{\sum_{j=0}^{N} w_j} \cdot (t_1 - t_0)$$

# How fairness is implemented: vruntime

- Virtual runtime (vruntime)
  - Charge each task a runtime proportional to $w_{base}$ and inversely proportional to its weight $w_i$
- Tasks are scheduled in order of increasing vruntime

$$V_{T_i}(t_1) = \frac{w_{base}}{w_i} \cdot (t_1 - t_0)$$

# VTimeScheduler

# What else can we do?

# Implement good schedulers

# Implement Schedulers

Typically not in Java

Implement Schedulers

Typically not in Java

Source: Marcus Biel

# CFS was built in a simpler time

Just one L3 cache

- Much smaller CPUs
- Topologies much more homogeneous
- Cores spaced further apart, migration cost typically high
- Power consumption and die area wasn't as important
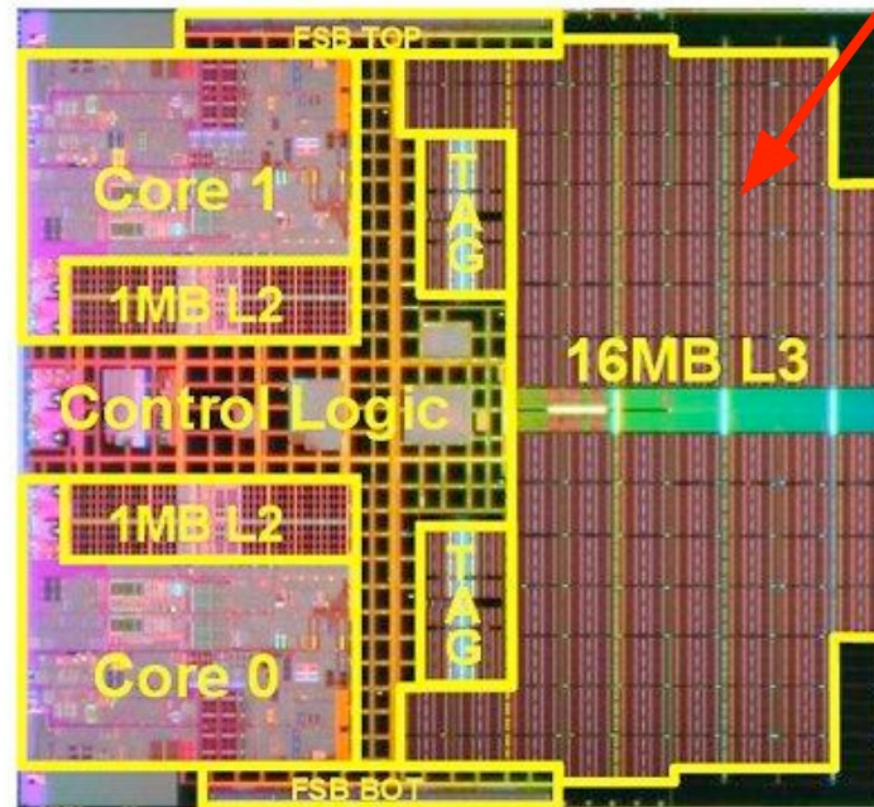- The fundamental assumptions behind heuristics may be easier to justify

Just two cores



FSB TOP
Core 1
1MB L2
TAG
Control Logic
16MB L3
1MB L2
TAG
Core 0
FSB BOT

Intel Xeon MP 71xx die

https://github.com/sched-ext/scx

# Reimplementing A Linux Rust Scheduler In eBPF Shows Very Promising Results

Written by Michael Larabel in Linux Kernel on 10 August 2024 at 03:27 PM EDT. 27 Comments

NVIDIA software engineer Andrea Righi has implemented his "scx_rustland" Linux Rust scheduler within eBPF for very promising performance results.

The bottleneck to the scx_rustland Rust-written scheduler has been the overhead in communication between kernel and user-space. To address this, he's implemented scx_rustland fully within eBPF and called the new creation scx_bpfland.

The scx_bpfland scheduler employs the same logic as scx_rustland but without the kernel/user-space communication overhead. Andrea has run some benchmarks and the new bpfland code is showing very promising results. PostgreSQL is as much as 30~39% faster, FFmpeg is several percent faster, nginx is around 8% faster, and more.

# scx_bpfland

# Gaming performance

- Frames per second (fps)
  - Primary metric for gaming performance
- Ideal fps for smooth gameplay
  - 30 fps: acceptable
  - 60 fps: fluid gaming experience
  - 120 fps: competitive gaming

# Experiments

# An erratic scheduler



https://lwn.net/SubscriberLink/1007689/922423e440f5e68a/

An erratic scheduler

Written in Java

Experiments

Having fun with sched-ext

## 4.3   Ensuring fair schedules

*Lol.*

All reasonable operating systems schedulers are *fair* —

# One that produces sound



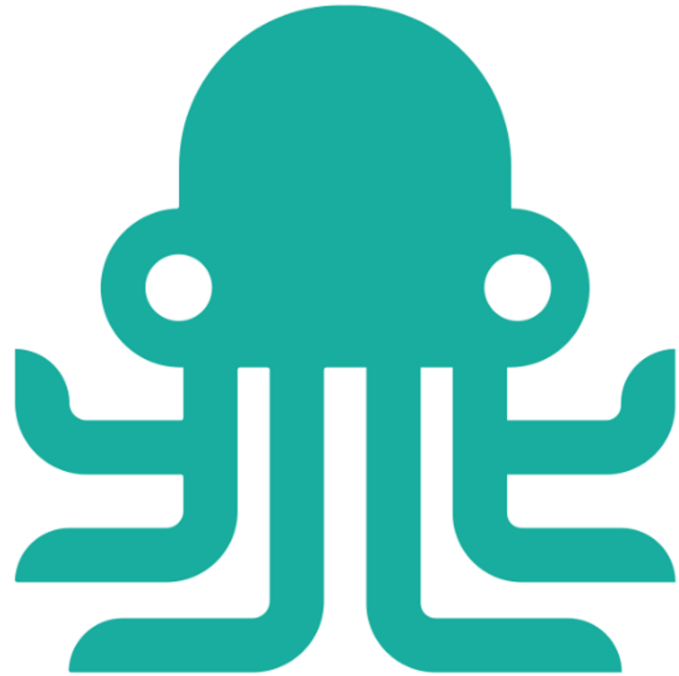https://github.com/parttimenerd/loudness-scheduler

# One that reacts to sound



https://github.com/parttimenerd/sound-of-scheduling

# TaskClicker



https://github.com/Mr-Pine/taskclicker

# Interactive,
# First Come , First Served
# Scheduler

# Interactive,
# First Clicked, First Served
# Scheduler

# The First Idle Game Scheduler

# TaskClicker

Failed after 32.236925450s

| | | | |
|---|---|---|---|
| Monitor Deflati 85667 -65ms | java 87658 -144ms | JS Watchdog 10535 -95ms | IPC I/O Parent 2396 -61ms |
| Isolated Web Co 13170 -38ms | systemd-udevd 539 -165ms | systemd-journal 495 -164ms | vesktop 70795 -124ms |
| Timer-0 85673 -27ms | Service Thread 85666 -26ms | Chrome_ChildIO T 70709 -135ms | Timer 78689 -136ms |

0 extra arms

0 eBees

Syscall balance: 313. Next upgrade at 2000

# Fin.

github.com/parttimenerd/hello-ebpf

**Johannes Bechberger**
**mostlynerdless.de**
OpenJDK Developer, SAP

**David Kiefer**
**mr-pine.de**
Student, KIT