Case for Transforming Parallel Run-times Into Operating System Kernels Paper Reading Group

> Kyle Hale Peter Dinda Presents: Maksym Planeta

> > 18.02.2016

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Table of Contents

Introduction

Evaluation

Development effort

Conclusion



Table of Contents

Introduction

Evaluation

Development effort

Conclusion

(4日) (個) (目) (目) (目) (の)

What is this project?

1. Northwestern University, Sandia Labs, Oak Ridge

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- 2. Part of Hobbes Project
- 3. They also develop Palacios

Proposes a microkernel



- Proposes a microkernel
- Uses hyperthreads in HPC context

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

- Proposes a microkernel
- Uses hyperthreads in HPC context

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 … のへで

Targets Xeon Phi

- Proposes a microkernel
- Uses hyperthreads in HPC context
- Targets Xeon Phi
- It cites L4 paper:
 - [40] J. Liedtke. On micro-kernel construction. In Proceedings of the 15th ACM Symposium on Operating Systems Principles (SOSP 1995), pages 237–250, Dec. 1995.

Idea

- $1.\ \mbox{HPC}$ app runs in user mode
- 2. Hardware available in kernel mode
- 3. When an HPC program runs in kernel mode:

(ロ)、(型)、(E)、(E)、 E) の(の)

3.1 All nice features are directly available

A Typical dialog with the kernel

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

ARE PROVIDED KERNEL ABSTRACTIONS THE RIGHT ONES?



ARE PROVIDED KERNEL ABSTRACTIONS THE RIGHT ONES?



RESTRICTED ACCESS TO HARDWARE



RESTRICTED ACCESS TO HARDWARE



Motivation

- 1. HPC app runs in user mode
- 2. Hardware available in kernel mode
- 3. When an HPC program runs in kernel mode:
 - 3.1 All nice features are directly available
 - 3.2 Kernel does not restrict the program with bad abstractions: For example, the run-time might need subset barriers, and be forced to build them out of mutexes.

Motivation

- 1. HPC app runs in user mode
- 2. Hardware available in kernel mode
- 3. When an HPC program runs in kernel mode:
 - 3.1 All nice features are directly available
 - 3.2 Kernel does not restrict the program with bad abstractions: For example, the run-time might need subset barriers, and be forced to build them out of mutexes.
 - 3.3 Kernel may waste resources for the features the application doesn't need:

For example, the run-time might not require coherence, but get it anyway.

Contributions

- 1. Criticize traditional architecture
- 2. Propose a new OS structure
- 3. Port some of the existing runtimes

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

```
Hybrid Run-Time (HRT)
```



(a) Current Model



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

3

► The runtime is the kernel, built within a kernel framework

◆□ ▶ < 圖 ▶ < 圖 ▶ < 圖 ▶ < 圖 • 의 Q @</p>

► The runtime is the kernel, built within a kernel framework

・ロト・日本・モト・モート ヨー うへで

Everything is kernel space

The runtime is the kernel, built within a kernel framework

(ロ)、(型)、(E)、(E)、 E) の(の)

- Everything is kernel space
- HRT has full access to the hardware

The runtime is the kernel, built within a kernel framework

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

- Everything is kernel space
- HRT has full access to the hardware
- HRT can pick its own abstractions

Aerokernel

User Mode



Figure 2: Structure of Nautilus.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Benefits

Better abstractions

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

- Noiseless
- Lightweight

Legacy support



(c) Hybrid Run-time Model Within a Hybrid Virtual Machine

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Table of Contents

Introduction

Evaluation

Development effort

Conclusion



Thread creation



Figure: Average, minimum, and maximum time to create a number of threads in sequence.

(日)、

æ

Thread creation



Figure: <u>Linux</u> from previous figure

・ロト ・四ト ・ヨト ・ヨト

Thread creation



Why bends? At (d) at 8 threads, (e) at 32, and (f) at 8. Bugs?

(日)、

æ

Thread creation (summary)

- iguio -. Noi uo

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

OS	Avg	Min	Max
Nautilus	16795	2907	44264
Linux	38456	34447	238866

Figure 3: Time to create a single thread measured in cycles.

Spinlock microbenchmark

OS	Execution time (s)	
Nautilus	13.72	
Linux	12.53	
OS	Avg. acquire/release time (cycles)	
Nautilus	59	
Linux	36	

Figure 5: Total time to acquire and release a spinlock 500 million times on Nautilus and Linux, and average time in cycles for an acquire/release pair.

Wake-up microbenchmark



Figure 6: Average event wakeup latency.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Circuit simulator benchmark



Figure 11: Run time of Legion circuit simulator versus core count. The baseline Nautilus version has higher performance at 62 cores than the Linux version.

Circuit simulator benchmark



Figure 12: Speedup of Legion (normalized to 2 Legion processors) circuit simulator running on Linux and Nautilus as a function of Legion processor (thread) count.

Circuit simulator benchmark



Figure 13: Speedup of Legion circuit simulator comparing the baseline Nautilus version and a Nautilus version that executes Legion tasks with interrupts off.

Table of Contents

Introduction

Evaluation

Development effort

Conclusion



Kernel development

The process of building Nautilus as minimal kernel layer with support for a complex, modern, many-core x86 machine took six person-months of effort on the part of seasoned OS/VMM kernel developers.

Language	SLOC
С	22697
C++	133
x86 Assembly	428
Scripting	706

Figure 9: Source lines of code for the Nautilus kernel.

Run-time support

Porting Legion:

- ▶ 43000 SLOC in C++
- Most of the work went into understating Legion
- Some code added to Nautilus

Language	SLOC
C++	133
С	636

Figure 10: Lines of code added to Nautilus to support Legion, NDPC, and NESL.

Four person-months to port

Also porting NESL and NDPC (related to each other).

Table of Contents

Introduction

Evaluation

Development effort

Conclusion



Conclusion

- A mikrokernel
- And a lightweight kernel
- Requires effort for porting
- Early stage of development

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ