# Tuning application performance with dynamically scalable parallelism

## Paper **#XX**

NN pages

## Abstract

In recent years, economic and physical realities have pushed both chip designers and application developers to rely increasingly on parallelism to satisfy performance demands [8]. However, there is no free lunch. Even after an application developer has wrestled with several target architectures, she still does not know what resources the application will get at runtime. Often, it is easier to over-provision units of execution and hope for the best. Yet even mild contention can drastically reduce overall system throughput and wall-clock execution times of individual applications. On the other hand, under-provisioning units of execution means leaving resources unused and falling short of performance potential. In this paper, we show that by selectively and dynamically changing application parallelism, we can avoid both pitfalls [in some way].

Stronger guarantees from scheduler to application (monitor actual resource utilization). Insights from application to scheduler (e.g., some information about app?).

## 1.  Questions

- How do developers typically provision number of threads? How do they determine the maximum?

## 2.  To do

- DS Paper outline

- MS AC Understand difference in performance for different levels of parallelism

- DS YW Understand Intel TBB

  - http://cas.ee.ic.ac.uk/people/dt10/teaching/2012/hpce/hpce-lec6-tbb-intro-handout-4up.pdf

  - https://drive.google.com/file/d/0B_10gtxnPV-_UFUyRGQ4MTQ4YTA/view?usp=sharing

- Design interface between TBB and coordinator

- Implement and integrate coordinator with TBB

- Integrate TBB into applications

  - http://wiki.cs.princeton.edu/index.php/PARSEC

- Evaluations

## 3.  Introduction

### What's the problem

- Parallel hardware + software hard to reason about

- Too many combinations to hand-tune

- Existing approach either over or under provision parallelism

### Why does it matter?

- Parallelism is here to stay both in hardware and software

- Challenges listed above make more important to have automated methods (we don't program in assembly anymore unless we have to)

### Who cares?

- Cloud providers / clusters pack many applications onto machines

- In large applications that take up entire machine, easier to create many threads

- In distributed systems, also easier since performance bottlenecks are elsewhere, but can see similar per-node improvements

### Contributions

- Observe that dynamically scaling parallelism can improve system performance (throughput, latency)

- User-level coordinator to determine appropriate level of parallelism for registered applications

- Thread library that can dynamically scale application parallelism
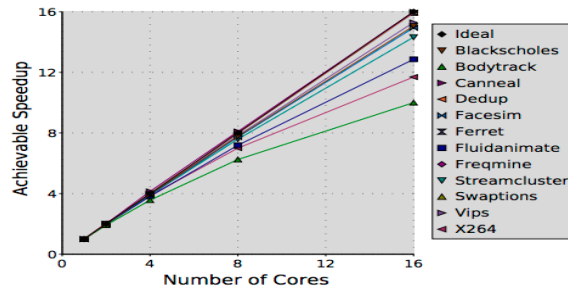
**Figure 1.** PLACEHOLDER: Wall-clock speed-up time. Upper bound for speedup of PARSEC workloads with input set simlarge based on instruction count. Limitations are caused by serial sections, growing parallelization overhead and redundant computations.
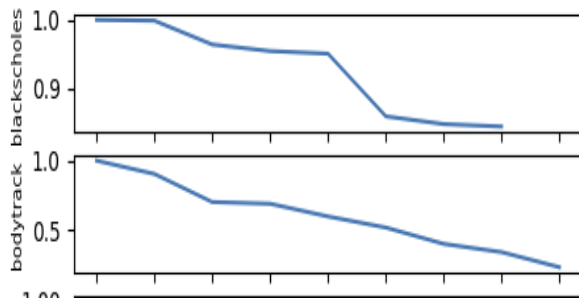


**Figure 2.** PLACEHOLDER: Although applications typically enjoy wall-clock speed-up times when they use more threads, often see corresponding increase in total CPU time [this graph might be confusing since we're talking about increase in total CPU time, but curve is going down]

## 4. Motivation for a dynamic parallelism

- Many cores are here to stay. Need more automated methods to take advantage of. [2]
- OS schedulers already quite complex. Must serve many masters; hard to make changes without harming certain workloads. Don't want to change the kernel if we can improve performance in user space [7]
- Applications are being increasingly run in containers / VMs on a single host
- Centralized schedulers in distributed systems struggle to scale [9]

## 5. Performance of parallel programs

Note: will cover PARSEC/Splash2 in more detail in later section [4] [3] [5]

### 5.1 Performance without contention

- speed-up curve (see Figure 1)
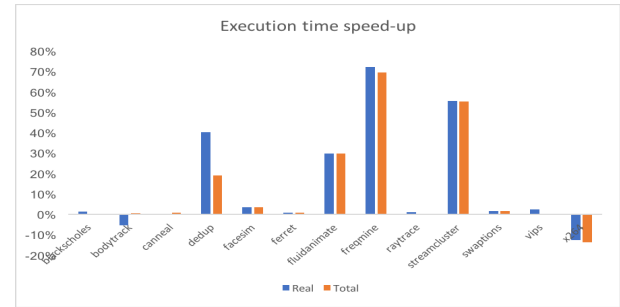- slow-down curve (see Figure 2)



**Figure 3.** PLACEHOLDER: How groups of applications (same or mixed (need this fig)) do under contention at varying # of threads

### 5.2 Performance under contention

- same application (see Figure 3)
- mixed group of applications (see Figure [])

### 5.3 Reasons for contention

- thread contention (manny runnable threads, but divided work up into small pieces; not overlapping well with i/o)
- lock contention (why would this be exacerbated under contention?)
- other resource contention

#### Classes of application behavior

- Big improvement from reducing parallelism (dedup, facesim, freqmine, streamcluster)
- Slight improvement from reducing parallelism
- Performance hit from reducing parallelism (x264)

### 5.4 Common strategies for dealing with unknown resources

TODO: maybe delete this section?

- Overprovision number of threads and hope for the best. But this can cause problems (e.g., lock or resource contention)

## 6. Determining contention and who should change parallelism

TODO: change section title
　　TODO: maybe this shouldn't be its own section

### 6.1 System metrics for determining contention

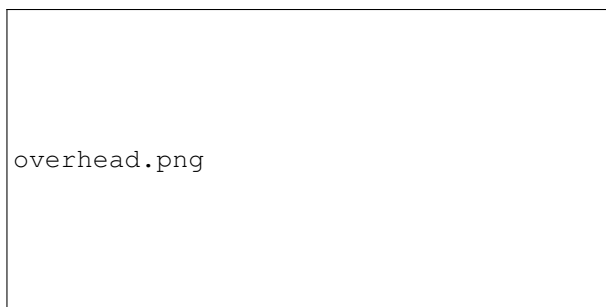- # of runnable threads
- CPU utilization

**Figure 4.** What is performance overhead of thread library over pthreads?

### 6.2 Application metrics for determining candidates for parallelism scaling

## 7. Design

Not sufficient to set parallelism just at application start, need to scale up and down. Machine load may change over time. Application requirements may change within a single run

### 7.1 Thread library

Similar to Intel Thread Building Blocks [10], but dynamically scale number of workers across applications, rather than just across workers

- Types of blocks: iterative / non-iterative parallelized operations
- DAG
- Unstructured

### 7.2 Driver application

Track CPU utilization by application (any other metrics?) Change target parallelism (how to decide, when to decide, how much to change by, which to change) Change CPU share or quota

### 7.3 Client applications

## 8. Benchmark applications

- TODO: do we want a single large application (e.g., provisioned for KNL, but running on much smaller VM for some reason)?
- TODO: What about an unstructured application (e.g., web server)?
- TODO: What about long-lived application (e.g., streaming)? We can probably simluate this via streamcluster

### 8.1 PARSEC

### 8.2 SPLASH-2

TODO: maybe leave this out given timing

## 9. Implementation

## 10. Evaluation

### 10.1 Workloads

Different classes of applications. Mixture of applications.
Should we have single large application?
How should we handle unstructured parallel programs (e.g., web server?)

### 10.2 Comparison with other thread libraries

OS (pthreads), OpenMP, Intel TBB

### 10.3 Varying guarantees (CPU share / quota)

### 10.4 Various machines

## 11. Related work

- Intel TBB, OpenMP, Cilk [6]
- Scheduler activations [1]
- Multikernel
- Exokernel
- Hierarchical schedulers
- Cache-aware / contention-aware / numa-aware schedulers
- Arachne
- Capriccio
- Wasted cores

## 12. Experience and future work

## Acknowledgements

## References

[1] ANDERSON, T. E., BERSHAD, B. N., LAZOWSKA, E. D., AND LEVY, H. M. Scheduler activations: Effective kernel support for the user-level management of parallelism. *ACM Transactions on Computer Systems (TOCS) 10*, 1 (1992), 53–79.

[2] BAUMANN, A., BARHAM, P., DAGAND, P.-E., HARRIS, T., ISAACS, R., PETER, S., ROSCOE, T., SCHÜPBACH, A., AND SINGHANIA, A. The multikernel: a new os architecture for scalable multicore systems. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles* (2009), ACM, pp. 29–44.

[3] BIENIA, C., KUMAR, S., SINGH, J. P., AND LI, K. The parsec benchmark suite: Characterization and architectural implications. In *Proceedings of the 17th international conference on Parallel architectures and compilation techniques* (2008), ACM, pp. 72–81.

[4] BIENIA, C., AND LI, K. *Benchmarking modern multiprocessors*. Princeton University, 2011.

[5] BIENIA, C., AND LI, K. Characteristics of workloads using the pipeline programming model. In *Computer Architecture* (2012), Springer, pp. 161–171.

[6] CONTRERAS, G., AND MARTONOSI, M. Characterizing and improving the performance of intel threading building blocks. In *Workload Characterization, 2008. IISWC 2008. IEEE International Symposium on* (2008), IEEE, pp. 57–66.

[7] LOZI, J.-P., LEPERS, B., FUNSTON, J., GAUD, F., QUÉMA, V., AND FEDOROVA, A. The linux scheduler: a decade of wasted cores. In *Proceedings of the Eleventh European Conference on Computer Systems* (2016), ACM, p. 1.

[8] MACK, C. A. Fifty years of moore's law. *IEEE Transactions on semiconductor manufacturing 24*, 2 (2011), 202–207.

[9] OUSTERHOUT, K., WENDELL, P., ZAHARIA, M., AND STOICA, I. Sparrow: distributed, low latency scheduling. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles* (2013), ACM, pp. 69–84.

[10] REINDERS, J. *Intel threading building blocks: outfitting C++ for multi-core processor parallelism.* " O'Reilly Media, Inc.", 2007.