

# 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 propose a mechanism to coordinate CPU usage across applications and show that by selectively and dynamically changing application parallelism, we can avoid both pitfalls.

## 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
- ✓ Understand Intel TBB
  - 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 context / what are the trends

- Increasing number of cores. Many-core architectures are here to stay. [2]
- As a result, application developers must think about parallelism to get performance
- Cloud becoming increasingly popular; run many containers / VMs on a single host

### What happens today

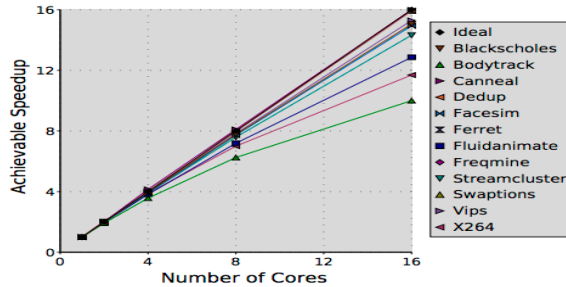
- Application: assume have full machine (over-provision) because not easy to know what resources we have.
- Application: different applications have different speed-up curves
- OS: schedule individual threads (Linux). Some control with cgroups
- OS: schedulers already quite complex. Must serve many masters; hard to make changes without harming certain workloads. Don't want to change the kernel if we can improve performance in user space [7]

### What's are the issues

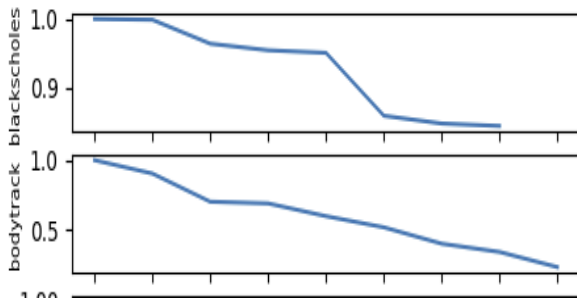
- Over-provisioning may low utilization on many core
- Also not taking advantage of best parallelism (slow-down in total CPU time)
- share and quota (cgroup) not enough; want to change parallelism
- Parallel hardware + software hard to reason about; too many combinations to hand-tune; need automated run-time methods

### Contributions

- Observe that dynamically scaling parallelism can improve system performance (throughput, latency)
- Mechanism for coordinating among applications (is this a contribution?)



**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]

- Thread library that can dynamically scale application parallelism
- User-level coordinator to determine appropriate level of parallelism for registered applications

## 4. Performance of parallel programs

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

### 4.1 Performance without contention

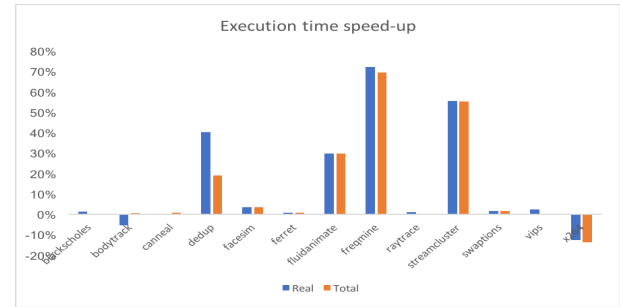
- speed-up curve (see Figure 1)
- slow-down curve (see Figure 2)

### 4.2 Performance under contention

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

### 4.3 Reasons for contention

- thread contention (manny runnable threads, but divided work up into small pieces; not overlapping well with i/o)



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

- 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)

## 4.4 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)

## 5. Determining contention and who should change parallelism

TODO: change section title

TODO: maybe this shouldn't be its own section

### 5.1 System metrics for determining contention

- # of runnable threads
- CPU utilization

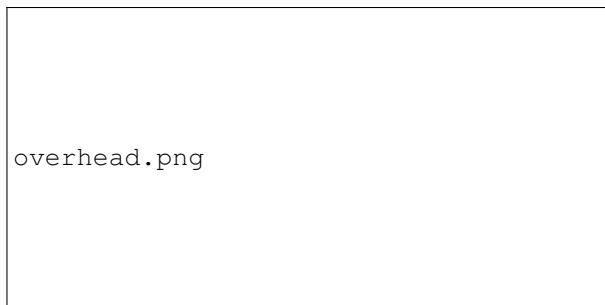
### 5.2 Application metrics for determining candidates for parallelism scaling

## 6. 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

### 6.1 Thread library

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



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

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

I/O and CPU pools treated differently in most libs

## 6.2 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

## 6.3 Client applications

## 7. 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 simulate this via streamcluster

## 7.1 PARSEC

## 7.2 SPLASH-2

TODO: maybe leave this out given timing

## 8. Implementation

## 9. Evaluation

### 9.1 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?)

### 9.2 Comparison with other thread libraries

OS (pthreads), OpenMP, Intel TBB

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

### 9.4 Various machines

## 10. Related work

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

## 11. 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.