

# Dynamic Core Processors

Kyle Daruwalla, David McNeil, and Ben Schmidt

*Rose-Hulman Institute of Technology*

**Abstract**—As both general purpose, desktop-grade processors and embedded processors mature, the distinguishing gap between them continues to diminish. Out of this narrowing market, a need for a versatile, low-power core emerges. Multicore embedded processors are just over the horizon. Dynamic core processors attempt harness these low-power cores to maximize both parallel throughput and minimize serial latency.

**Keywords**—*IEEEtran, journal, L<sup>A</sup>T<sub>E</sub>X, paper, template.*

## I. INTRODUCTION AND MOTIVATION

Mobile computing devices have grown to require processors that can support a dynamic workload. The typical serial tasks such as voice compression, speech transcoding (for communications), and image compression must perform well. However, newer mobile devices need to render complex graphics and run advanced background scheduling, all while maintaining serial performance and power. The obvious answer might be GPUs, but they are known to be power-hungry chips. Ideally, if most of the serial and parallel tasks could be performed on a single chip, the GPU would only need to be used to perform high TLP tasks like displaying graphics. A smaller GPU workload means a lower power GPU. Thus, power consumption can be kept minimal, while performance gains are still attained.

Dynamic core processors attempt solve precisely these issues. By adding an interconnection network, some additional hardware, and control logic, the symmetric multicore embedded processors can be reconfigured into a modern, super-scalar, out-of-order processor. Thus, by identifying serial and parallel program sections, the processor can be dynamically changed to perform efficiently.

Obviously, creating such a processor, though possible, is not a trivial task. So, we attempt to analyze the theoretical performance gains of a dynamic core as described above. The following work will show whether a dynamic core processor performs better than a symmetric multicore processor and a modern, super-scalar, out-of-order processor. Furthermore, for a fixed length program, we determine how often the program must switch from serial to parallel processing before a dynamic core processor realizes performance gains. Finally, we determine the maximum number of cores in a dynamic processor before critical path length negatively impacts serial performance.

## II. PRIOR WORK

Summarize relevant related work with appropriate citations.

## III. METHODOLOGY

### A. Models

In order to simulate a dynamic core processor without actually implementing one, the team identified two basic

Benchmark	Description	Parameters
ocean	placeholder	placeholder
fft	placeholder	placeholder
lu	placeholder	placeholder
radix	placeholder	placeholder

TABLE I: The SPLASH-2 benchmarks used

models - an baseline symmetric core and a modern, out-of-order core. By combining multiple baseline cores into a single multiprocessor, we could create a parallel model (i.e. a processor biased towards parallel programs). The single modern, out-of-order core makes our serial model (i.e. a processor biased towards serial programs). Thus, the third model in our study, the dynamic model, is a combination of the serial and parallel models.

Based on this methodology, we needed a simulator capable of switching between CPU types. Previous work in this area had been done using the gem5 simulator. Capable of switching between CPU types, gem5 was an ideal choice.

### B. Tools

The gem5 simulator [1] is a combination of M5, a simulation framework with support for multiple ISAs and CPU models, and GEMS, a memory simulation system in two parts - Ruby and Opal. As a result, gem5 is a robust simulator with support for five ISAs, ARM, ALPHA, MIPS, Power, SPARC, and x86, and four CPU models:

- AtomicSimple is a minimal single IPC CPU model
- TimingSimple is similar to AtomicSimple but also simulates the timing of memory references,
- InOrder is a pipelined, in-order CPU
- O3 is a pipelined, out-of-order CPU model

gem5 has two modes of operation - system emulation mode and full system mode. System emulation mode does not model the OS or peripheral devices but solely simulates the specified benchmark. Full system mode, on the other hand, uses an actual OS kernel and mounts a Linux disk image. Essentially, gem5 full system mode is capable of booting a full OS and presenting the user with a Linux command prompt.

Initially, gem5's many high level customizations suggested that it would be an effective simulator for a dynamic processor model. However, while gem5 may boast many impressive features, we found that many of these features are not fully supported or difficult to configure.

Benchmark	Description	Parameters
jpeg	placeholder	placeholder
epic	placeholder	placeholder
gsm	placeholder	placeholder
adpcm	placeholder	placeholder

TABLE II: The MediaBench II benchmarks used

For simulation of a dynamic processor, we needed two types of benchmarks. One which would be representative of highly parallelized workloads and one reflecting serial execution. The parallel benchmark we used was SPLASH-2 [2]. SPLASH-2 is a suite of benchmarks intended to evaluate the performance of multiprocessor systems. For a serial benchmark suite, we chose to use MediaBench II [3]. A suite of compression and decompression multimedia algorithms. The programs selected from each benchmark were intended to accurately reflect typical mobile device usage. Tables I and II detail the exact benchmarks used.

### C. Procedure

When we initially began the project we planned on using gem5 SE mode for simplicity. However, due to the parallel nature of SPLASH-2 a multi-threading library is required. Because SE mode does not emulate a full operating system, there was no OS-level threading library such as pthreads. As a result, we finally decided to use full system mode. This provided us with access to Linux’s full threading library allowing us to simulate the multi-threaded applications. Using full system mode had the additional benefit of providing extremely accurate results because the benchmarks are actually running in a Linux operating system, similar to our application space.

At this point, it was necessary to determine what CPU types and ISAs would be tested in order to simulate a dynamic processor. We initially planned on using either x86 or ARM as the underlying ISA. However, we soon found that only AtomicSimple CPU model was supported by gem5 full system mode for these ISAs. gem5 provides much more complete support for the ALPHA ISA. ALPHA is a 64-bit RISC ISA representative of a simply architected general processor. Using this ISA as our building block we were able to develop benchmarks to represent a serial and a parallel machine. Our serial system used the O3 CPU model to simulate a modern out of order processor. Our parallel system was comprised multiple TimingSimple CPU models. Originally, we had intended to use the InOrder model, however it was not available. Even though the TimingSimple CPU type was not a standard in-order pipeline, an array of the single cycle cores out performed the out-of-order processor on most parallel benchmarks. Thus, we felt the loss in accuracy was negligible.

In order to efficiently run tests, we developed scripts to boot a Linux kernel in FS mode and run through each of our benchmarks. At the completion of each benchmark the

timing statistics would be written out to file. A script was then developed capable of parsing the output file and produce a CSV file containing the statistics for each benchmark. The Linux boot process became an advantage in our setup. Spanning 2.4 trillion instructions, we could be confident that all our caches would be warmed up after the boot process was complete. Furthermore, gem5’s built-in stats-dump tool allowed us to run serial and parallel benchmarks one after the other without shutting the system down. This meant that the cache state before each benchmark was exactly the cache state at the end of the previous benchmark; thus, the initial cache miss latencies were embedded into our timing results.

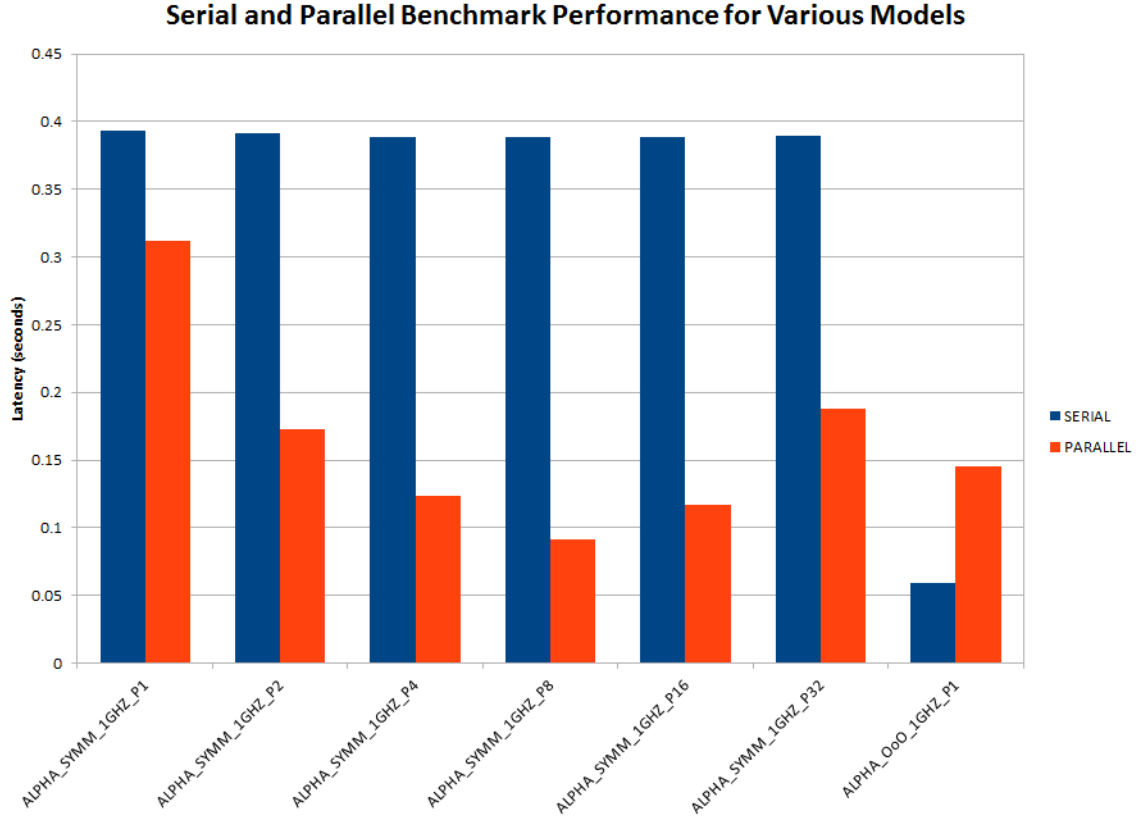
Once we had the timing results for the serial and parallel models, we synthesized the results together to get the timing performance for a dynamic core. Of course, there is some latency associated with switching from a serial configuration to a dynamic configuration. Though we did not implement a dynamic core, we did theorize potential methods for a datapath. In particular, we were drawn to the idea of using muxes to run external inputs into the execution units of several cores. Based on this, it seemed appropriate that an advanced interconnect network would be required. Most multicore and SoC devices use a crossbar network, but the crossbar is hardly scalable. Instead, many bleed-edge devices use network-on-chip (NoC) structures that are based off of crossbar interconnects, but optimized for scalability and power dissipation. Thus, we selected a switching latency of 0.991 ns based on research in the field of NoCs [5].

Furthermore, in order to study the problems surrounding a dynamic processor, we tackled to optimization issues. First, we chose a single serial program and a single parallel program to make up a single “chunk” of mixed code. We then look run time for 32 chunks put together. Then, the size of the serial and parallel portion in a chunk is doubled, so that each chunk is twice as large as it was before. However, we now only use 16 chunks in the test, so that the overall number of instructions remains constant. Repeating this process, we are able to get results for the peak number of switches for maximum performance between serial and parallel portions within a fixed code segment. Secondly, we recognized that as the number of cores in the parallel model increases, the critical path of the serial model increases. Thus, the frequency should scale down for the serial model. Based on simple geometry (Pythagorean’s theorem), we assume that the frequency will scale on the order of  $1/\sqrt{n}$ . Thus, we obtained results for the serial model with scaled frequencies to see if there is a drop-off in performance.

## IV. RESULTS

### A. Model Accuracy

First, we needed to verify that our serial and parallel models were accurate. As expected, the serial model out performed the parallel model at serial tasks. Additionally, the parallel model performed the same for a serial task regardless of the number of cores; this is expected, since only a single core can effectively be used. Finally, the parallel model performed better than the serial model for parallel programs. However, it



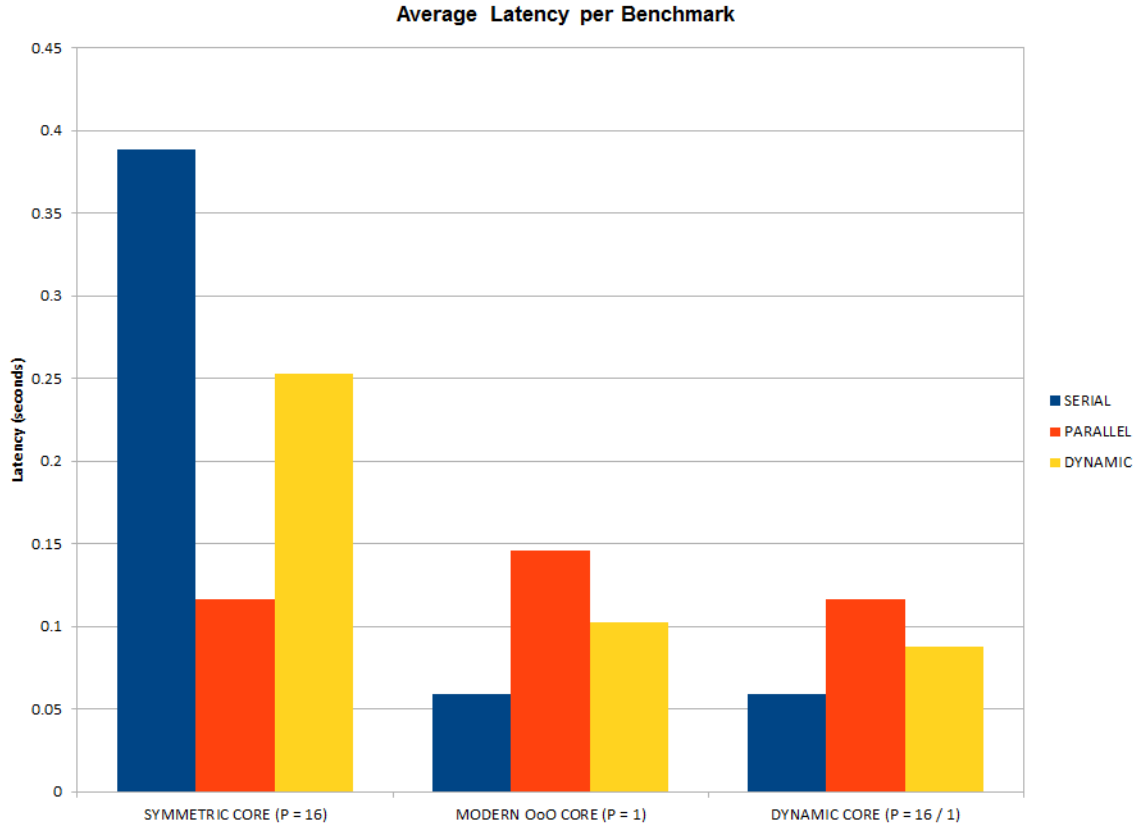
**Fig. 1:** The serial and parallel models perform appropriately and accurately for different benchmarks

is worth noting that as the number of cores increased passed eight, the parallel performance began to degrade. At 32 cores, the performance was even worse than the serial model. Even so, we feel that these models are accurate for our theoretical study. See Figure 1 for these results.

#### B. Is Dynamic Better?

Ultimately, we wanted to answer the question "Is Dynamic Better?". At this stage, "better" is simply evaluated in terms of performance (we will go into more detail later). As can be seen in Figure 2, a dynamic core performs better in terms of average latency over the dynamic benchmark. Based on these results, the dynamic core has about a 4x improvement over a single baseline core machine (a baseline core being the cores used to make up the parallel model), a 2.9x improvement over a sixteen core symmetric multicore machine, and 1.2x improvement over a modern, out-of-order processor. Note that we assumed the TimingSimple model in gem5 was sufficient to be used as our baseline core. This conclusion is corroborated by the results in Figure 1; yet, if we had used a RISC pipelined architecture instead of a single cycle processor, the parallel model performance would have been even better (note how small the difference is between the parallel model and the serial model for the parallel benchmark in Figure 2). So, if anything, our dynamic core would have even greater performance gains.

However, as we briefly mentioned before, "better" is more than just performance. Especially on a mobile device, "better" is determined by efficiency. While we did not quantitatively determine the power consumption, we can qualitatively say a dynamic core will be more power efficient by nature. Since the symmetric multicore is made up of small, low power cores, and these same cores make up the basis for a dynamic core, a dynamic core is bound to be power efficient. There is, of course, overhead associated with a dynamic core. But this overhead should be no more than the overhead involved in a modern, out-of-order processor. We can analytically estimate this power overhead. Based on research at MIT, the maximum total power for a flip-flop or latch in a high performance system is 350 uW. Given that we are simulating an eight issue-width modern, out-of-order processor, we can expect about eight functional units (functional units scale with issue width). Since each function unit requires two inputs and one output (32-bits each) to be muxed, each functional unit will have  $(3muxes)((32latches) + (32flip - flops)) = 192latches/flip - flops$ . So, the average power overhead is  $0.350(8)(192) = 537.6$  mW, which is minimal compared to the power consumption of a modern, out-of-order processor. So, speculatively (with some evidence), we can say that a dynamic core should consume less power than a modern, out-



**Fig. 2:** A dynamic core performs better than a 16 symmetric multicore machine and a single modern out-of-order machine

of-order processor, while still performing better.

### C. Peak Number of Switches

When discussing a processor that performs context switching, it is useful to mention granularity. In other words, how often should a processor be switching between serial and parallel mode. As Figure 3 shows, the finer grained the context switching, the lower the performance gains. Do not let the scale of the graph be misleading, the problem has purposefully been scaled up. So, our smallest building block between a context switch is a full program long (i.e. a full FFT or voice compression algorithm). Thus, the degradation is not as drastic numerically. However, logically, the trend should hold as the granularity becomes finer. Since the only factor degrading the performance is a linear constant for the interconnect latency, it follows that very fine grained switching will lower performance gains visibly.

### D. Frequency Scaling

Though increasing the cores in a dynamic processor will result in performance gains due to the parallel mode running faster, the serial portion will begin to suffer. A large number of cores means a greater area, which means a longer critical path for the serial model. Based off Pythagorean's

theorem, we hypothesized that as the number of cores increased, the frequency would scale on the order of  $1/\sqrt{n}$ . After re-running our original tests, we determined the new dynamic core latencies with a scaled frequency serial model. The speedup relative to a 1 GHz modern, out-of-order core can be seen in Figure 4. There is a peak performance gain at four cores, before the speedup decays. So, while a dynamic core is better in an ideal sense, when practical constraints are applied to it, performance begins to degrade.

### E. Practicality

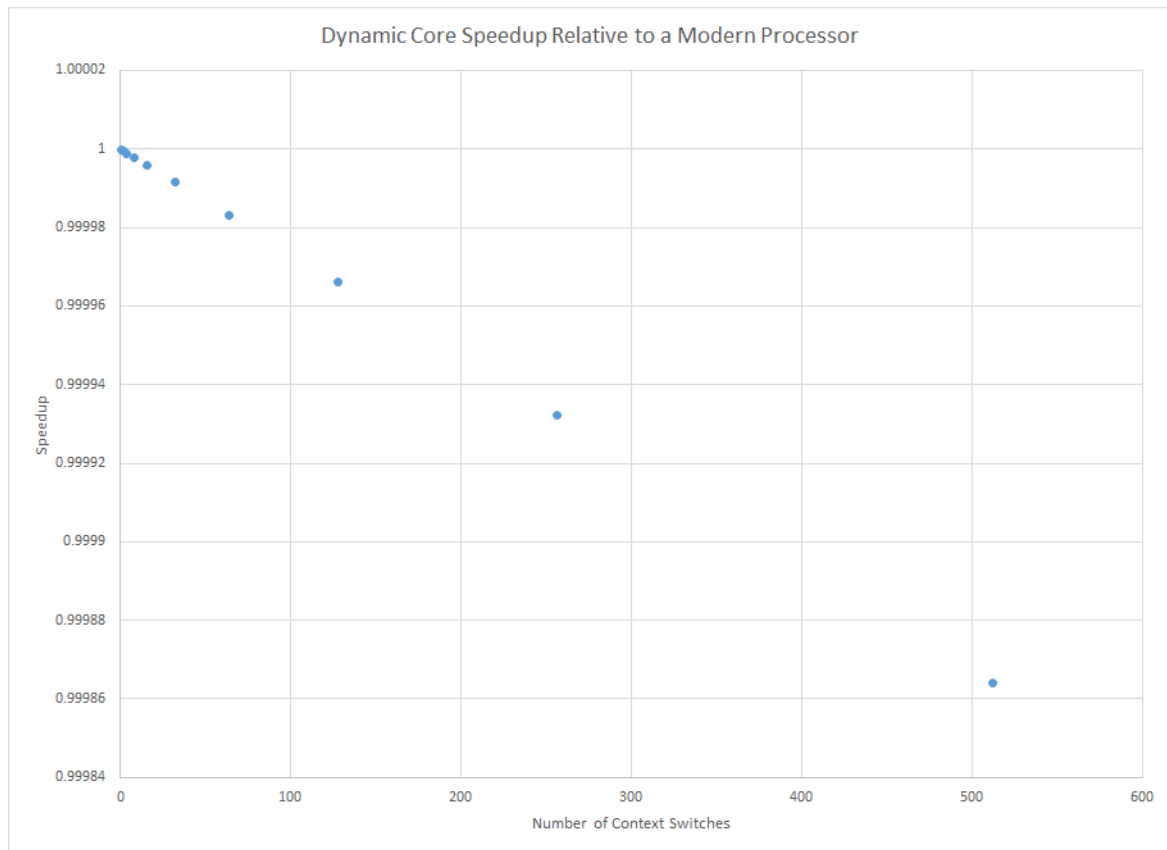
In addition to our theoretical models, we chose to compare our dynamic core to modern, out-of-order, dual core processor, which is typical for the average device today. The speedup was 0.75x, so a decrease in performance.

## V. CONCLUSION AND FUTURE WORK

Summarize your findings, identify remaining open problems or issues, and suggest the next steps for this project.

## VI. STATEMENT OF WORK

The project report must also include a statement of work that identifies the contributions of each individual on the team.



**Fig. 3:** *Fine granularity switching degrades performance gains*

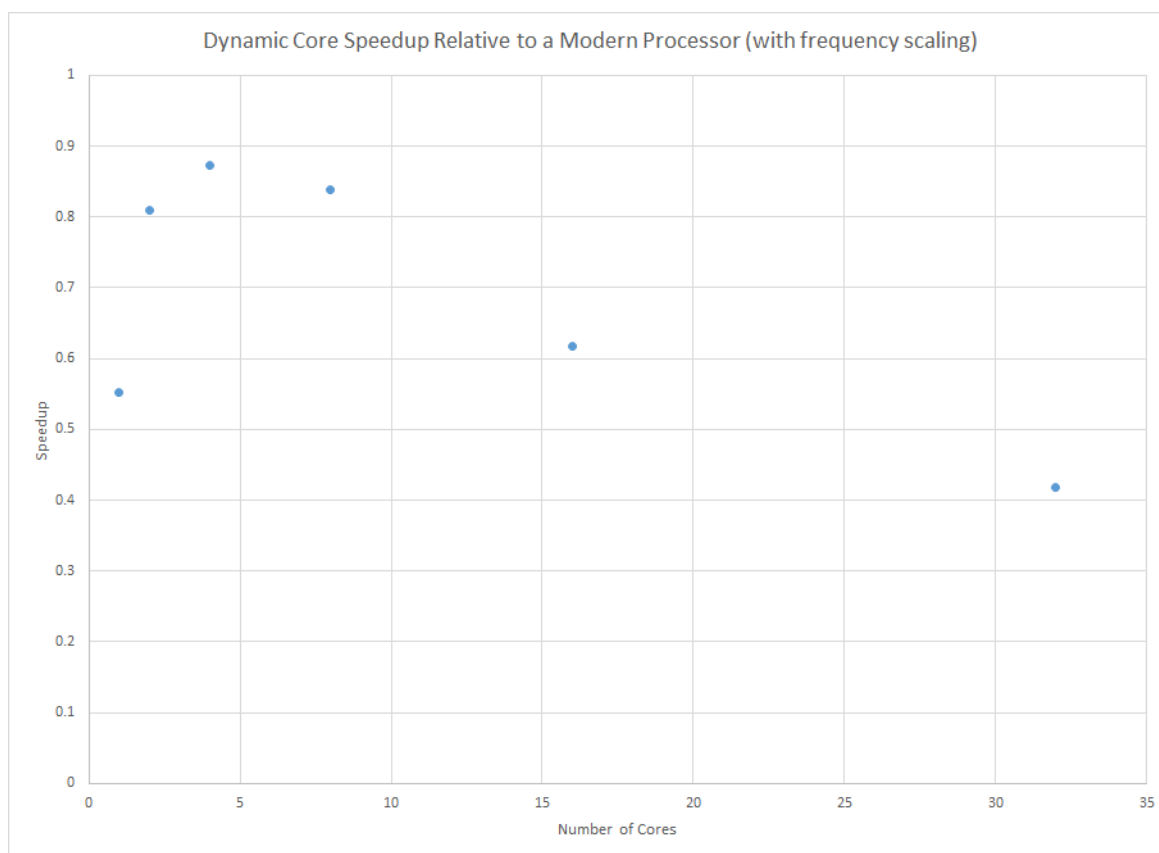
The statement of work must reflect a team consensus and must be signed by all team members. I recommend that you structure this statement as a table with a row for each project milestone, a column for each team member, and the percentage contribution of each team member to each milestone in the entries in the table. Please do not give a vanilla 50/50 in every column.

#### ACKNOWLEDGMENT

The authors would like to thank...

#### REFERENCES

- [1] N. Binkert, B. Beckmann, G. Black, S. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. Hill, D. Wood. "The gem5 Simulator" [Online] Available: [http://research.cs.wisc.edu/multifacet/papers/can11\\_gem5.pdf](http://research.cs.wisc.edu/multifacet/papers/can11_gem5.pdf)
- [2] Not sure how to site this. Available: <http://www.capsl.udel.edu/splash/index.html>
- [3] Not sure how to site this. Available: <http://euler.slu.edu/~fritts/mediabench/>
- [4] Mark D. Hill, Michael R. Marty. "Amdahls Law in the Multicore Era" [Online] Available: [http://moodle.rose-hulman.edu/pluginfile.php/245005/mod\\_resource/content/0/Amdahl\\_Multicore%20%28Hill%29.pdf](http://moodle.rose-hulman.edu/pluginfile.php/245005/mod_resource/content/0/Amdahl_Multicore%20%28Hill%29.pdf)
- [5] Sung-Joon Lee, Jaeha Kim. "A 256-Radix Crossbar Switch Using Mux-Matrix-Mux Folded-Clos Topology" [Online] Available: [http://www.jsts.org/html/journal/journal\\_files/2014/12/Year2014Volume14\\_06\\_10.pdf](http://www.jsts.org/html/journal/journal_files/2014/12/Year2014Volume14_06_10.pdf)
- [6] Vladimir Stojanovic, Vojin Oklobdzija, Raminder Bajwa. "Comparative Analysis of Latches and Flip-Flops for High-Performance Systems" [Online] Available: [http://www.rle.mit.edu/isg/documents/Stojanovic\\_ICCD98.pdf](http://www.rle.mit.edu/isg/documents/Stojanovic_ICCD98.pdf)



**Fig. 4:** Increasing the number of cores scales down the serial model frequency, resulting in a peak performance at four cores