

Research Dissertation

Andrew Wright [40006161]

May 4, 2013

Abstract

Here's an abstract

1 Introduction

This document will investigate how best to analyse the performance of task data flow and work-first based runtime systems for parallel programs taking into account data locality.

These runtime systems use a task data flow model to record and manage dependencies during parallel execution. This approach allows the programmer to concentrate on application development rather than the low level parallelisation techniques. This is useful because processing different tasks on different cores can be tricky, especially if there are implicit dependencies between tasks. They also provide an effective and flexible way for the scheduler to maximise the number of tasks that can be run at one time across the multiple processors. Using this method, we can rival highly tuned libraries with minimal tuning effort [13]. The other aspect of these schedulers is their provably good work first principle. Which will be discussed in further

detail in section 1.1.2. The scheduler analysed is the first scheduler of this kind and is called 'Swan' [15].

To give a proficient introduction to the problem area I will discuss why it is useful to combine task dependencies with the work-first principle by looking at schedulers which focus on one method in particular and then present the unified solution which this research will analyse. It is important to understand some of the finer points of these runtime systems as they will shape how the analysis tool is built. I will then discuss the existing methods for analysing parallel runtime systems focusing on methods which were considered in this research. The rest of the paper will go onto set out my approach and results.

1.1 The Problem Area - Parallel Scheduling

1.1.1 Task Based Parallel Programming

An implementation of task based parallel programming is *SMPSs* [13]. It attempts to provide a way for the programmer to avoid the arbitrary concerns of data dependencies and

parallelisation as much as possible allowing them to focus on the problem at hand. To this effect, it is implemented as an extension of the C programming language with additional pragmas to identify functions as candidates to be run in parallel.

SMPSs uses the idea of a task dependency DAG (directed acyclic graph). Where each vertex is a task and the edges represent dependences between tasks. SMPSs uses what we call a non-strict DAG which means that one task can have a dependency on another which isn't on the same ancestral tree. We will discuss this concept further in Section 1.1.3.

SMPSs adopts a “work-stealing” approach whereby if a thread has no more tasks to execute, either on its own queue or the system queue, it will work on tasks from another thread's queue in FIFO order. Work-stealing is also implemented in Swan and Cilk as we will see.

A potential downside of SMPSs, compared to work-first schedulers (discussed later), is that it favours the critical path over work. Its master thread creates all dependent tasks for a function leaving the scheduler with an excess of idle tasks. There is a theory set out by the developers of Cilk [6] which proves that, based on an acceptable degree of parallel slack, it is actually more efficient to adopt a “work-first” approach. This means sacrificing critical path time i.e. more scheduling overheads, in order to reduce the amount of time it takes to complete a run on one processor. I will talk more about how this is achieved in the next section.

1.1.2 Work-First Scheduling

Cilk is where the idea of work-first scheduling originated and an implementation of it can be found in Cilk-5 [6]. Cilk is a parallel extension of C. The parallelism is introduced particularly through spawn and sync statements. Spawns allow a procedure to be run in parallel with its caller and sync requires the procedure to stall until the completion of all spawned procedures. The aspects of Cilk which we are particularly interested in, to provide a background on Swan, are the work-stealing and work-first properties of Cilk. Cilk uses a provably good work-stealing algorithm developed in Cilk-1 [2, 3, 10].

This “work-stealing” algorithm is implemented by giving each processor a ready deque (a double ended queue) data structure of call frames (a stack of related frames). Call frames can be inserted from the bottom and removed from either end. The processor uses the ready deque as a call stack, pushing and popping from the bottom. Processors which steal from this processor remove a call frame from the top. The reason that it steals from the top of the victim's (the processor from whom the call frame is being stolen) deque is to avoid contention and increase the probability of data locality.

The main development of Cilk-5 is a focus on the “work-first” approach to scheduling. We can characterise a Cilk computation by two quantities: its work, which is the total time needed to execute the computation serially, and its critical-path length, which is its execution time on an infinite number of processors. As its name suggests the “work-first”

principle is:

“Minimise the scheduling overhead borne by the work of a computation. Specifically, move overheads out of the work and onto the critical path.”

In order for the work first principle to be effective we must assume that there is ample “parallel slackness”. This means that the number of processors required to achieve full parallelism for the program is sufficiently greater than the number of available processors. For a mathematical definition of ‘work’, ‘parallelism’ and ‘critical path’ please refer to section 2.2. The reason that Cilk needs this slackness is to ensure that the extra expense of less efficient scheduling algorithm can be made up for by the efficiency of the work done per processor.

As previously mentioned Cilk enforces a fully strict DAG which means that direct dependencies cannot be created from one task to another which isn’t on the same ancestor tree. This is unlike Swan where the task data flow structure allows dependencies between tasks on differing ancestor trees. See Figure.1 for a comparison. Because Cilk uses a strict DAG it means that it is best suited for applications that can adopt a ‘divide and conquer’ approach to problem solving and is most effective when there is regular parallelism throughout the program. It isn’t possible to write some programs, which could benefit from parallel execution, in a this way for example pipelines.

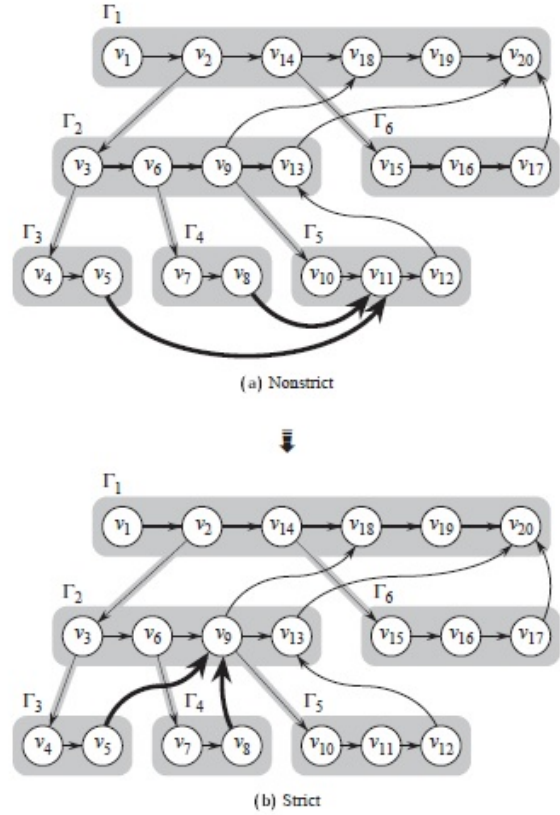


Figure 1: Comparison between strict and non-strict DAGs. Accredited to PhD Thesis by Robert D. Blumofe [2]

1.1.3 A Unified Scheduler

The scheduler which this research is focused on is the Swan scheduler. Swan combines work-first scheduling (the approach which Cilk adopts) with dependency-aware scheduling (as seen in SMPSSs).

Swan's language is similar to Cilk and the dependency aspect of the scheduler is realised at the object level. The concept of a versioned object is introduced here. These objects are passed to tasks as 'indep', 'outdep' and 'inoutdep' arguments depending on their memory access type. These are defined if a given argument has a memory access mode of: input, output or input/output for a task. If a task has no dependency information associated with it, it is run as an unconditional spawn and is executed as it would be with Cilk.

The swan scheduler allows nesting of fork/join parallelism and task graph parallelism. As well as task graphs which are arbitrarily nested. Nested task graphs require care because if the same object is being referenced by both a parent and a child the versioned object meta-data must track the correct dependencies. Swan handles this by creating new counters in the child's object so that dependencies are tracked independently.

1.2 Scheduler Analysis

The aim of this research is to analyse parallel schedulers so I will now describe the most relevant work which has been done in this area. There are several analysing tools for parallel systems and the majority follow a similar

process which debugs and analyses a post-mortem parallel program. This process is:

1. *Data collection* : Software is used to collect information about the execution of the parallel program (hardware can be used too).
2. *Data analysis* : Once the stream of events is collected, the analysis step can be started to order, filter, calculate some statistics or simulate a parallel machine.
3. *Display* : Visualise the result provided by the analysis step.

I will look at each of the scheduler analysing tools in terms of these three areas.

1.2.1 General Parallel Program Analysis

There have been tools developed which can analyse the execution trace of a given parallel scheduler. Paraver is such a tool, it visualises parallel execution and displays it in a customisable way. Paraver looks at the actual execution time as a basis of its analysis. It uses various windows to display information to the user all of which is customisable so the user can view the collection of windows which most suits their purpose. These update as the user traces through the program, representing what is happening through various views. The ability to filter and to visualise execution is useful however a lot of effort is required from the user in order to get a meaningful metric to help them improve their application. Paraver is also only interested in the

display aspect of the analysis process. The data collection and exporting have to be applied to the program before Paraver is able to display it from a simulation file. This collection is very expensive and collecting a trace file can slow a runtime system down by an order of magnitude - leaving the user with an inaccurate measurement.

1.2.2 Work First Analysis

Cilkview is an analysing tool for Cilk which is useful to this research because Swan extends Cilk and benefits from much of the same analysis.

Cilkview first collects a program's parallelism information during a serial execution of the program code. It uses meta data in Cilk++ to identify the parallel control constructs in the executing application which allows it to construct a data dependency DAG. As this meta data is built into the Cilk++ binary it imposes no overhead on performance in a normal production environment. Instead of capturing exact timing measurements Cilkview uses instruction count as a surrogate of time. They feel that parallelism estimates accurate to within a binary order of magnitude suffice to diagnose problems of inadequate parallelism.

Cilkview focuses on three areas when analysing the results from *work*, *span* and *burdened span*. I will look at how each of these are calculated in turn and discuss why they are important.

Work is the total amount of time spent in all the strands. Let T_P be the fastest possible execution time of the application on P pro-

cessors. Since the work corresponds to the execution time on 1 processor, we denote it by T_1 . One reason that work is an important measure is that it provides a lower bound on P-processor execution time:

$$T_P \geq T_1/P$$

This is at least how much time it takes for P processors to execute the total amount of work done. Because theoretically P processors can do P instructions at one time. Using Cilk T_1/T_P is the most *speedup* that can be gained. Speedup is how much faster a program can run on a given environment this can be intuitively written as a program which is run on a machine with 5 processors can run 5 times as fast. This level of speedup is known as linear speed up.

The second area which Cilkview analyses is span. Span is also called the critical path and is the longest amount of time which is taken for a program to execute over any path of dependencies in the DAG. The critical path is the theoretical fastest time which a program could be run on an unlimited amount of processors. Due this it is denoted as T_∞ this also provides a bound on P processor execution time:

$$T_P \geq T_\infty$$

Cilkview uses the work and span to calculate the parallelism of a program T_1/T_∞ . The parallelism is the average amount of work along each step of the critical path. From parallelism we can work out what the maximum number of processors are on which linear speed-up can be achieved for a particular program. Linear speed-up cannot be

achieved on any number of processors greater than the parallelism. This information allows Cilkview to determine how scalable a program is.

The third aspect which Cilkview analyses is the burdened span. Although the execution time is reduced by tasks being run on various processors the overhead and book-keeping which is introduced to schedule these executions needs to be taken into account. This is what the burdened DAG does. In Cilkview a constant called the *span coefficient* δ is multiplied by the span to take into account the work stealing overhead which includes the book keeping and the cost of cache misses in order to migrate the stolen task's working set.

$$T_P \leq T_1/P + \delta T_\infty$$

The burdened DAG places a burden on each continuation and return edge of the DAG. Cilkview computes the *burdened span* by finding the longest path in the burdened DAG.

If we let T_1 be the work of an application program, and let \hat{T}_∞ be its burdened span. Then, a work-stealing scheduler running on P processors can execute the application in expected time:

$$T_P \leq T_1/P + 2\delta\hat{T}_\infty$$

Cilkview uses this burdened model to give a more accurate estimated lower bound on speedup. This is, a work-stealing scheduler running on P processors achieves speedup at least:

$$\frac{T_1}{T_P} \geq \frac{T_1}{T_1/P + 2\delta\hat{T}_\infty}$$

A detailed proof of this can be found in Cilkview's research paper [8].

Cilkview displays the following information to the user: amount of work done, the span, the burdened span, the parallelism, the burdened parallelism, number of spawns, number of syncs and the average maximal strand (which is how much work is down in between scheduler operations). Cilkview then gives estimated speedups for 2 - 32 processors. Cilkview also uses graphs to display this data in a meaningful way.

1.3 Other Metrics

There are other metrics which have been used to measure the performance of parallel applications, examples of these are: True Zeroing[11], Gprof[7] and Quartz NPT[1]. An overall study and comparison of each method has been carried out [9] and suggests that critical path analysis, as discussed in section 1.2.2, is most effective.

2 Investigation Method

As alluded to above the most useful metrics which currently exist to analyse parallel programs revolve around the program's span also known as its *Critical Duration*.

2.1 Critical Duration Analysis

The critical duration is the shortest possible time takes for a given parallel program to execute. This is often represented as T_∞ as it is the time taken on an infinite number of processors. In this research the task graph will be analysed at run time and the critical duration established.

This information allows us to calculate the degree of parallelism which, as seen in Section 1.2.2, will determine how scalable a program is. Scalability is a vital measurement in regards to parallel programming as the primary goal of parallel programming is to allow performance to increase as the number of threads increase.

In order to get a value for parallelism, and therefore scalability, another metric is required. This is the work done. Work is defined as the total amount of time spent in all strands. This can be denoted as T_1 that is the time taken when the program is run on a single thread because a serial execution will carry out all of the work sequentially.

2.1.1 Measurement of Time

There are two options to consider when measuring time. The first is instruction count which is the method that Cilkview [8] adopts. They argue that because the advantages of adding another processor would only be realised if there was a large difference between parallelism and actual speed-up there is no need for a precise measurement of time.

The other option is to take the actual time which is the method that this research will

use. The benefit of this is that time is an easily relatable to metric and can be validated externally. Time was captured at a millisecond level to give an accurate result. This was chosen over cycles in order to increase portability and to ensure that changes brought about by CPU power saving options where the frequency is reduced wouldn't affect the measurements. Also `gettimeofday()` from the `sys/time.h` library was used over `RDTSC` [?] due to issues with CPU timers not being synchronised over multiple cores.

A potential negative of this method of timing is the fact that it is it will take into account the entire host environment rather than the individual process. So if the user is running the program on a machine which is heavily loaded with running tasks the time measurements could be slower than on the same machine when it running only the program. This could affect the level of parallelism to be expected. It would still be accurate however in terms of wall clock time which is what we are concerned with here.

2.1.2 Non-object Dependent Critical Duration Analysis

The critical path algorithm devised in this research can be broken into two parts. Firstly the work first aspect, whereby object dependencies are ignored and secondly the object aware aspect. We will now discuss the details of these aspects.

The algorithm for obtaining the critical duration at run time without object dependencies can be seen in Figure 2. The parent's work is calculated on the current strand by

taking the current time away from the last start time which will give us the parent's current duration, we then add that to the current amount of work done by the parent to get the current total work done. The parent's critical path is then updated. All work that parent has done on the current strand will get directly added into their critical path. The spawned/called task is assigned the parent's current critical duration. We will see the reason for this when objects are introduced.

When a task is finished and returns to it's parent the duration, work done and critical duration are calculated. The calculation of the critical duration takes into account the child's critical duration also. If a task has spawned children who have had a longer critical duration than the task itself then their critical duration should become it's critical duration. Once this has been determined the task's duration should be added to the critical duration because, like the parent tasks, the duration of a given task will always be a part of it's critical duration.

2.2

References

- [1] T.E. Anderson and E.D. Lazowska. *Quartz: A tool for tuning parallel program performance*, volume 18. ACM, 1990.
- [2] R.D. Blumofe. *Executing multithreaded programs efficiently*. PhD thesis, Massachusetts Institute of Technology, 1995.

when each task is ready to execute:

```
calculate parent's duration
calculate parent's work done [p.
    wd + duration]
calculate parent's critical path
    [p.cp + duration]
```

```
assign task's critical path =
    parent's critical path
```

when a task is finished executing:

```
calculate task's duration
calculate task's work done [t.wd
    + duration]
calculate task's critical
duration:
    t.cd = MAX(t.
        child_critical_duration
        , t.cd)
    t.cd = t.cd + duration
```

```
calculate parent's critical
duration:
    if the task has been
        spawned:
        p.cd += t.cd
    else break
calculate parent's child critical
duration:
    p.ccd = MAX(p.ccd, t.
        critical_duration)
```

```
increase parent's work done by
    task's work done
```

when a task is synced

```
parent's critical duration is
    calculated [MAX(p.cd, p.
        child_critical_duration)]
```

Figure 2: Algorithm describing the critical duration as captured without object dependencies

- [3] R.D. Blumofe, C.F. Joerg, B.C. Kuszmaul, C.E. Leiserson, K.H. Randall, and Y. Zhou. *Cilk: An efficient multithreaded runtime system*, volume 30. ACM, 1995.
- [4] C. Ding and T. Chilimbi. A composable model for analyzing locality of multi-threaded programs. Technical report, Technical Report MSR-TR-2009-107, Microsoft Research, 2009.
- [5] M. Frigo, P. Halpern, C.E. Leiserson, and S. Lewin-Berlin. Reducers and other cilk++ hyperobjects. In *Proceedings of the twenty-first annual symposium on Parallelism in algorithms and architectures*, pages 79–90. ACM, 2009.
- [6] M. Frigo, C.E. Leiserson, and K.H. Randall. The implementation of the cilk-5 multithreaded language. *ACM Sigplan Notices*, 33(5):212–223, 1998.
- [7] S.L. Graham, P.B. Kessler, and M.K. McKusick. Gprof: A call graph execution profiler. *ACM Sigplan Notices*, 17(6):120–126, 1982.
- [8] Y. He, C.E. Leiserson, and W.M. Leiserson. The cilkview scalability analyzer. In *Proceedings of the 22nd ACM symposium on Parallelism in algorithms and architectures*, pages 145–156. ACM, 2010.
- [9] J.K. Hollingsworth and B.P. Miller. Parallel program performance metrics: a comparison and validation. In *Proceedings of the 1992 ACM/IEEE conference on Supercomputing*, pages 4–13. IEEE Computer Society Press, 1992.
- [10] C.F. Joerg. *The cilk system for parallel multithreaded computing*. PhD thesis, Citeseer, 1996.
- [11] B.P. Miller, M. Clark, J. Hollingsworth, S. Kierstead, S.S. Lim, and T. Torzewski. Ips-2: The second generation of a parallel program measurement system. *Parallel and Distributed Systems, IEEE Transactions on*, 1(2):206–217, 1990.
- [12] Stephen L. Olivier, Bronis R. de Supinski, Martin Schultz, and Jan F. Prins. Characterizing and mitigating work time inflation in task parallel programs. 2012.
- [13] J. Pérez, R. Badia, and J. Labarta. A dependency-aware task-based programming environment for multi-core architectures. pages 142–151, 2008.
- [14] Andreas Sandberg, David Black-Schaffer, and Erik Hagersten. Efficient techniques for predicting cache sharing and throughput. In *Proceedings of the 21st international conference on Parallel architectures and compilation techniques*, PACT ’12, pages 305–314, New York, NY, USA, 2012. ACM.
- [15] H. Vandierendonck, G. Tzenakis, and D.S. Nikolopoulos. A unified scheduler for recursive and task dataflow parallelism. pages 1–11, oct. 2011.