

# *Oracle-Guided Scheduling for Controlling Granularity in Implicitly Parallel Languages*

Umut A. Acar

*Carnegie Mellon University, Pittsburgh, PA, USA  
Inria, Paris, France  
(e-mail: umut@cs.cmu.edu)*

Arthur Charguéraud

*Inria, Université Paris-Saclay, France  
LRI, CNRS & Univ. Paris-Sud, Université Paris-Saclay, France  
(e-mail: arthur.chargueraud@inria.fr)*

Mike Rainey

*Inria, Paris, France  
(e-mail: mike.rainey@inria.fr)*

---

## **Abstract**

A classic problem in parallel computing is determining whether to execute a thread in parallel or sequentially. If small threads are executed in parallel, the overheads due to thread creation can overwhelm the benefits of parallelism, resulting in suboptimal efficiency and performance. If large threads are executed sequentially, processors may spin idle, resulting again in suboptimal efficiency and performance. This “granularity problem” is especially important in implicitly parallel languages, where the programmer expresses all potential for parallelism, leaving it to the system to exploit parallelism by creating threads as necessary. Although this problem has been identified as an important problem, it is not well understood—broadly applicable solutions remain elusive.

In this paper, we propose techniques for automatically controlling granularity in implicitly parallel programming languages to achieve parallel efficiency and performance. To this end, we first extend a classic result, Brent’s theorem (a.k.a. the work-time principle) to include thread-creation overheads. Using a cost semantics for a general-purpose language in the style of lambda calculus with parallel tuples, we then present a precise accounting of thread-creation overheads and bound their impact on efficiency and performance. To reduce such overheads, we propose an oracle-guided semantics by using estimates of the sizes of parallel threads. We show that, if the oracle provides accurate estimates in constant time, then the oracle-guided semantics reduces the thread-creation overheads for a reasonably large class of parallel computations.

We describe how to approximate the oracle-guided semantics in practice by combining static and dynamic techniques. We require the programmer to provide the asymptotic complexity cost for each parallel thread and use runtime profiling to determine hardware-specific constant factors. We present an implementation the proposed approach as an extension of the Manticore compiler for Parallel ML. Our empirical evaluation shows that our techniques can reduce thread-creation overheads, leading to good efficiency and performance.

## 1 Introduction

*Explicit parallel programming* provides full control over parallel resources by offering primitives for creating, scheduling, and synchronizing parallel threads. As a result, the programmer can, at least in principle, write efficient parallel programs by performing a careful cost-benefit analysis to determine which threads should be executed in parallel and under what conditions. This approach, however, requires reasoning about low-level execution details, such the effects of scheduling, data races or concurrent effects, which is known to be notoriously hard. In addition, such low-level programming can lead to over-engineering of software and harm portability. For example, the resulting code may perform well on a particular hardware setting but not on others.

The complexities of parallel programming with explicit languages have motivated interest in *implicitly parallel languages*, such as Multilisp (33), NESL (17), Cilk (19), Parallel Haskell (21; 43), Parallel ML in Manticore (27; 28) and in MultiMLton (49). The key distinction between implicitly parallel programs and explicitly parallel programs is that implicitly parallel programs accept a sequential (serial) semantics, whereas explicitly parallel programs may not. Implicitly parallel programming thus dramatically simplifies writing parallel programs. Indeed languages for implicit parallelism enable the programmer to express parallelism at a high level via language constructs, e.g., parallel sequences, parallel arrays, parallel tuples, and delegate the job of parallel execution to the compiler and the run-time system.

As an implicit parallel program executes, it exposes opportunities for parallelism (as indicated by the parallel constructs). The language run-time system exploits the potential for parallelism by creating parallel threads as needed and mapping them to the processors or cores for fast execution. To achieve efficiency, implicit programming languages rely on a scheduler for distributing parallel threads among the processors to perform load balancing. Various scheduling techniques and practical schedulers have been developed, including work-stealing schedulers (e.g., (19)), and depth-first-search schedulers (16).

Experience with implicitly parallel programs shows that one of the most important decisions that any implicit parallel language must make is determining whether or not to exploit an opportunity for parallelism by creating a parallel thread or forgo the opportunity by falling back to sequential execution. On the one hand, creating a thread for each parallelism opportunity can lead to poor efficiency, because the cost of creating and managing parallel threads can be very high. On the other hand, foregoing a parallelism opportunity can lead to suboptimal performance because of the lost opportunity for parallelization of a large chunk of work. Therefore, solving this dilemma requires finding just the right “sweet spot”, where no more than necessary threads are created to achieve the best performance. Solving this problem, referred to as the *granularity problem*, is important because the overhead of parallel threads matters: since the speedups achievable via parallel computation is bounded by the number of processors, often a small constant factor, any increase in the overhead, however small, can negatively affect performance.

### MANUAL GC

Granularity problem has been studied extensively. One approach, developed in early work in functional programming is to use cost prediction to determine the threshold values. One approach uses list size to determine the threshold values (35) but assume linear work

complexity for every parallel operation. Another approach uses the depth and height of the recursion tree (54; 41) to determine the threshold values. Depth and height are not good predictors of computation cost, however. Lopez *et al.* take this approach as well, but in the context of logic programming (38).

One approach to the granularity problem is to focus on reducing the cost of creating a thread, rather than limiting the number of threads created. This approach is taken by implementations of work stealing with lazy thread creation (39; 25; 26; 30; 45; 34; 47; 52). In lazy thread creation, the work stealing scheduler is implemented so as to avoid, in the common case, the major scheduling costs, in particular, those of inter-processor communication. But, in even the most efficient lazy thread creation, there is still a non-negligible scheduling cost for each implicit thread.

Lazy Binary Splitting (LBS) is an improvement to lazy thread creation that applies to parallel loops (52). The crucial optimization comes from extending the representation of a thread so that multiple loop iterations can be packed into a single thread. This representation enables the scheduler to both avoid creating closures and executing deque operations for most iterations. A limitation of LBS is that it addresses only parallel loops whose iteration space is over integers. Lazy Tree Splitting (LTS) generalizes LBS to handle parallel aggregate operations that produce and consume trees, such as map and reduce (12). LTS is limited, however, by the fact that it requires a special cursor data structure to be defined for each tree data structure.

### 1.1 Amortizing Per-Thread Costs

Aharoni *et al.* study the granularity problem in the setting of distributed computing (8), where the crucial issue is how to minimize the cost of inter-processor communication. In their setting, the granularity problem is modeled as a staging problem, in which there are two stages. The first stage consists of a set of processor-local thread pools and the second stage consists of a global thread pool. Moving a thread to the global thread pool requires inter-processor communication. The crucial decision is how often each processor should promote threads from its local thread pool to the global thread pool. We consider a different model of staging in which there is one stage for parallel evaluation and one for sequential evaluation.

The approach proposed by Aharoni *et al.* is based on an online algorithm called CG. In this approach, it is assumed that the cost of moving a thread to the global thread pool is an integer constant, called  $g$ . The basic idea is to use amortization to reduce the scheduling total cost of moving threads to the global thread pool. In particular, for each thread that is moved to the global thread pool, CG ensures that there are at least  $g + 1$  threads added to the local thread pool. Narlikar describes a similar approach based on an algorithm called DFDeque (40). Just as with work stealing, even though the scheduler can avoid the communication costs in the common case, the scheduler still has to pay a non-negligible cost for each implicit thread.

### **1.2 Flattening and Fusion**

Flattening is a well-known program transformation for nested parallel languages (18). Implementations of flattening include NESL (17) and Data Parallel Haskell (42). Flattening transforms the program into a form that maps well onto SIMD architectures. Flattened programs are typically much simpler to schedule at run time than nested programs, because much of the schedule is predetermined by the flattening (50). Controlling the granularity of such programs is correspondingly much simpler than in general. A limitation of existing flattening is that certain classes of programs generated by the translation suffer from space inefficiency (16), as a consequence of the transformation making changes to data structures defined in the program. Our transformation involves no such changes.

The NESL (17) and Data Parallel Haskell (42) compilers implement fusion transformation in order to increase granularity. Fusion transforms the program to eliminate redundant synchronization points and intermediate arrays. Although fusion reduces scheduling costs by combining adjacent parallel loops, it is not relevant to controlling granularity within loops. As such, fusion is orthogonal to our oracle based approach.

Even though the granularity problem is broadly accepted to be an important problem, there has not been mu

it is a poorly understood one. Theoretical analyses often ignore thread-creation overheads, offering no significant clues about how such overheads may affect efficiency. Practical implementations often focus on reducing thread-creation overheads, instead of controlling granularity so that fewer threads can be created without harming parallelism. As a result, practitioners try to solve the granularity problem by estimating the amount of work that would be sufficiently large to justify parallel execution. More specifically, programmers try to determine the input sizes at which threads become too small to amortize the cost of parallel thread creation and sequentialize such threads. Since the running time of a thread depends on the hardware, the programmer must make the best decision they can by taking into account the specifics of the hardware. This manual granularity control is thus bound to yield code whose performance is not portable (i.e., platform specific) (52).

In this paper, we propose theoretical and practical techniques for solving the granularity problem in implicitly parallel-programming languages. First, we present theorems that take into account thread-creation overheads and characterize their impact on parallel run time. Our theorems show such impact to be significant (Section 4). We then consider a cost semantics for a calculus that extends the lambda calculus with parallel tuples (Section 5) and propose a technique for controlling granularity based on an oracle (Section 6). We show that if the oracle can be implemented efficiently and accurately, it can be used to improve efficiency for a relatively large class of computations.

Based on this result, we propose a practical realization of the oracle that uses complexity functions defined by the user to approximate accurately the work of the computations involved (Section 7). We then describe a compilation scheme for implementing the oracle-guided semantics based on such complexity functions (Section 8). We present an evaluation of the proposed approach on a subset of ML extended with parallel tuples and complexity annotations (Section 9). We conclude with a discussion of the limitations of our implementation of the oracle based on complexity functions (Section 11).

## 2 Overview

We present a high-level overview of our the techniques proposed in this paper.

Brent’s theorem (20), which leads to the *work-time scheduling principle*, characterizes arguably the most important property of parallel programs: that they can be scheduled to execute efficiently with multiple processors—within a factor of two of the optimal. More precisely, Brent shows that a level-by-level schedule can execute a parallel computation with  $w$  work and  $s$  span in no more than  $w/P + s$  steps on  $P$  processors. This theorem generalizes to any greedy schedule, that is, any schedule that never leaves processors idle when there is work to do. Prior research shows that greedy schedules can be computed by online scheduling algorithms such as the work-stealing algorithm (e.g., (19)) under modest assumptions.

The work-time scheduling principle ignores an important factor: the overhead for creating threads or parallelism, which is assumed to be zero. To understand the impact of thread-creation overheads, we start with this fundamental theorem and generalize it to take the overheads into account (Section 4). Specifically, we consider the standard directed-acyclic-graph (DAG) mode for parallel computations and show that a computation with *total work*  $\mathbb{W}$  and *total span*  $\mathbb{S}$ , where both include the thread-creation overheads, can be executed in no more than  $\mathbb{W}/P + \mathbb{S}$  steps. What is somewhat special about thread-creation overheads are that they are not divisible. Nevertheless, the generalized theorem follows by a modest extension of an existing proof.

Having established the contribution of thread-creation overheads to parallel run time, we then move on to the problem of determining precisely the overheads of thread creation in implicitly parallel programs. To this end, we consider a lambda calculus with parallel tuples and present a cost-semantics for evaluating expression of this language (Section 5). The *cost semantics* yield raw work/span and total work/span of each evaluated expression. Using this cost semantics, we show that thread-creation overheads can be significant: a multiplicative factor times the raw work. When applied to the generalized Brent’s theorem, this result implies that such multiplicative increases in work affect the parallel run-time directly. To reduce thread-creation overheads, we propose an alternative *oracle-guided semantics* that captures a known principle for avoiding the thread-creation overheads: for a thread, create a parallel thread and evaluate it in parallel only if the thread is sufficiently large, i.e., greater than some constant  $\kappa$ . We show that the oracle semantics can decrease the overheads of thread-creation by any (desired) constant factor  $\kappa$ , but only at the cost of increasing the total span by a similar factor. This result suggests that in practice some care will be needed to select  $\kappa$ , because otherwise it can reduce an important quantity called parallel slackness (53).

The bounds with the oracle-guided semantics suggest that we can reduce the thread-creation overheads significantly, if we can realize the semantics in practice. Such a realization is impossible, unfortunately, because it requires the ability to determine a priori the running time of a thread and do so without incurring other overheads. We show, however, that a realistic oracle that can give constant-factor approximations to the thread running times can still result in similar reductions in the overheads for a reasonably broad class of computations (Section 6.2). We also show that, unless care is taken, the realistic oracle can actually further increase the overheads, due to the direct cost of evaluating the oracle. This

outcome, i.e., that attempts at controlling the granularity can actually backfire and slow down the program further, is an interesting outcome of our analysis. For a broad class of computations, including many recursive, divide-and-conquer computations, we show that this worst case can be avoided.

While the oracle-guided semantics is effective in controlling the cost of thread creation without detrimentally harming parallelism, it is technically impossible to realize in practice because it requires predicting the work (sequential run-time) of computations. As we describe in Section 7, however, the work to be performed by a given thread can be approximated by using a combination of static and dynamic information. Specifically, we describe an approximation technique that relies on an estimator that uses asymptotic cost functions (asymptotic complex bounds) and judicious use of run-time profiling techniques to estimate actual run-times accurately and efficiently. This approach combines data from the asymptotic complexity bounds and the profiling techniques to approximate the work to be performed by a given thread. In this work, we only consider simple recursive functions for which the execution time is proportional to the value obtained by evaluating the asymptotic complexity expression. We refer to Section 11 for discussion of more general patterns.

We present a prototype implementation of the proposed approach (Section 8) by extending the ML to support parallel tuples, and complexity functions, and translating programs written in this extended language to the PML (Parallel ML) language (28). Although our implementation requires the programmer to enter the complexity information, these could also be inferred in some cases via static analysis (e.g., (36) and references thereof). We extend the Manticore compiler for PML to support oracle-guided scheduling and use it to compile generated PML programs. Our experiments (Section 9) show that our oracle implementation can reduce the overheads of a single processor parallel execution to 3 and 13 percent of the sequential time. When using 16 processors, we achieve 7- to 15-fold speedups on our benchmark machine.

### 3 Terminology

In this paper, we consider implicitly parallel programs, where a program executes by employing two kinds of threads: system-level threads and user-level threads. An execution creates one *system-level thread* per processor and usually pins the system-level thread to that processor. It then dynamically creates lighter-weight, *user-level threads* and maps them to the system-level threads. A key property of implicitly parallel programs is that the number of user-level threads created during an execution can be very large relative to the number of system-level threads. For example, an execution may create millions of user-level threads mapped to 10 system-level threads, which are then mapped to 10 processors or cores by the system. For the purpose of succinctness, throughout the paper, we refer to “user-level threads” simply as “threads” and to “system-level threads” as “processes.”<sup>1</sup>

<sup>1</sup> User-level threads are sometimes called “strands”; unfortunately, the same term is sometimes used to refer to hardware threads.

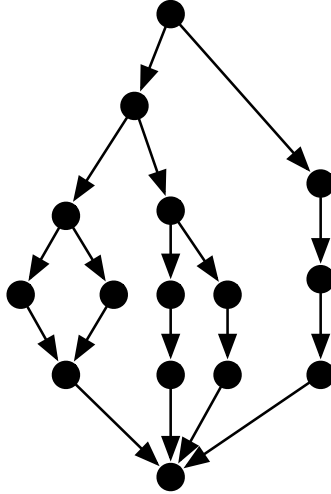


Fig. 1. An example computation dag.

#### 4 Generalizing Brent’s Theorem

We represent a parallel computation with a directed acyclic graph (dag), called *computation dag*. Vertices in the dag represent atomic computations, or operations. For simplicity, we refer to each vertex as an *operation*. Edges between operations represent precedence relations, in the sense that an edge from  $u$  to  $v$  indicates that the execution of  $u$  must be completed before the execution of  $v$  can start. Every computation dag includes a *source* operation and a *sink* operation, representing the starting and the end points of the computation, respectively. All operations of a computation dag are reachable from the source, and the sink is reachable from all operations. An example computation dag appears in Figure 1.

Traditionally, every operation in the dag is considered to take a single unit of time and given a weight of 1. In this setting, we can define the standard notion of work and span, which we here call *raw work* and *raw span*. The *raw work* of a computation graph is equal to the total number of operations that it contains. The *raw span* of the computation dag is equal to the total number of operations along the longest path. Brent proved the following bound, which we recall since our aim is to generalize it.

**Theorem 4.1 (Brent’s theorem)** *Let  $G$  be a computation dag with  $w$  raw work and  $s$  raw span. A level-by-level or depth-wise schedule of  $G$  requires  $O(\frac{w}{P} + s)$  time on  $P$  processors.*

*Proof*

Consider scheduling the operations *level-by-level*, that is, starting with operations at depth 1, then executing operations at depth 2, and so on, ending with operations at depth  $s$ . Let  $w_i$  denote the number of operations at depth  $i$  in the dag. These operations can be scheduled on  $P$  processors in  $\lceil \frac{w_i}{P} \rceil$  rounds, each lasting exactly one unit of time. Summing up over all the possible depths, we deduce that the total execution time is bounded by:

$$\sum_{i=1}^s \lceil \frac{w_i}{P} \rceil \leq \sum_{i=1}^s \left( \frac{w_i}{P} + 1 \right) \leq \frac{\sum_{i=1}^s w_i}{P} + s \leq \frac{w}{P} + s.$$

□

Observe that the optimal execution time cannot be less than  $\frac{w}{p}$ , which corresponds to having all processors busy at all time, and that it cannot be less than  $s$ , which corresponds to the length of a critical path. Therefore, Brent's bound, expressed as  $\frac{w}{p} + s$ , is never more than a factor 2 away from the optimal.

Brent's theorem does not take into account the overheads associated with thread creation. Our goal is to refine the model and generalize Brent's theorem to take thread-creation costs into account. To that end, we assign to each operation with out-degree two or greater a weight of  $1 + \tau$  instead of just 1. We then define the *total work* as the sum of the weights of all the operations in this revised computation dag. Similarly, we define the *total span* as the maximum weight of a path from the source to the sink in the revised dag.

Theorem 4.2 generalizes Brent's theorem for weighted dags and arbitrary greedy schedules (rather than level-by-level schedules). The proof is a relatively straightforward generalization of that of Arora et al (9).

**Theorem 4.2 (Greedy-scheduling of weighted dags)** *Let  $G$  be a computation dag with  $\mathbb{W}$  total work and  $\mathbb{S}$  total span. Any greedy  $P$  processors schedule of  $G$  takes  $O(\frac{\mathbb{W}}{P} + \mathbb{S})$  steps.*

*Proof*

Consider any greedy execution of the dag  $G$ . At each execution step, each processor places a token in the *work bucket* if it is busy at this step, otherwise it places a token in the *idle bucket*. The work bucket contains exactly  $\mathbb{W}$  tokens at the end of the execution. Let  $I$  denote the number of tokens contained in the idle bucket at the end of the execution, and let  $T$  denote the total number of steps in the execution. Because a total  $TP$  tokens are created during the  $T$  execution steps,  $TP = \mathbb{W} + I$ . In order to establish the result  $T \leq \frac{\mathbb{W}}{P} + \mathbb{S}$ , it thus suffices to establish the inequality  $I \leq P\mathbb{S}$ .

Consider a given time step. If all processors execute an operation at that step then the idle bucket receives zero tokens. Otherwise, a number of processors are idle. In this case, the idle bucket receives between one and  $P - 1$  tokens. If one or more processors are idle, then the idle processors cannot find a ready operation to execute, because otherwise this would not be a greedy schedule. This means that at this step, all the ready operations (operations whose ancestor have been executed) are executed. Therefore, at such a time step, the span of the sub-dag induced by the operations that are not yet executed decreases by one. Therefore, there are at most  $\mathbb{S}$  such steps and thus  $(P - 1)\mathbb{S}$  tokens in the idle bucket. □

## 5 Source Language and Cost Semantics

To give an accurate account of cost of thread creation, and to specify later our compilation strategy, we consider a source language in the style of  $\lambda$ -calculus and present a dynamic cost semantics for it. The semantics and the costs are parameterized by  $\tau$  and  $\phi$  that represent the cost of creating a parallel thread and the cost of consulting an external oracle.



$$\begin{aligned}
v & ::= x \mid \mathbf{n} \mid (v, v) \mid \mathbf{inl} \ v \mid \mathbf{inr} \ v \mid \mathbf{fun} \ f.x.e \\
e & ::= v \mid \mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2 \mid (v \ v) \mid \mathbf{fst} \ v \mid \mathbf{snd} \ v \mid \\
& \quad \mathbf{case} \ v \ \mathbf{of} \ \{\mathbf{inl} \ x.e, \mathbf{inr} \ x.e\} \mid (e, e) \mid (|e, e|)
\end{aligned}$$

Fig. 2. Abstract syntax of the source language

### 5.1 Syntax

The source language includes recursive functions, pairs, sum types, and parallel tuples. Parallel tuples enable expressing computations or *branches* that can be performed in parallel, similar to the fork-join or nested data parallel computations. More precisely, in  $(|e_1, e_2|)$ , we refer to  $e_1$  and  $e_2$  as (parallel) branches. We note that although we only consider parallel tuples of arity two, our results generalize to fixed (constant) arity trivially.

To streamline the presentation, we assume programs to be in A-normal form, with the exception of pairs and parallel pairs, which we treat symmetrically because our compilation strategy involves translating parallel pairs to sequential pairs. Figure 2 illustrates the abstract syntax of the source language. We note that, even though the presentation is only concerned with a purely-functional language, it is easy to add local mutable state (i.e. mutable memory cells that are not subject to concurrent accesses); in this case, however, they contribute no additional insight and thus are omitted for clarity.

### 5.2 Dynamic Cost Semantics

We define a dynamic semantics where parallel tuples are evaluated selectively either in parallel or sequentially, as determined by their relative size compared with some constant  $\kappa$ , called the cutoff value and such that  $\kappa \geq 1$ . To model this behavior, we present an evaluation semantics that is parameterized by an identifier that determines the *mode* of execution, i.e., sequential or not. For the purpose of comparison, we also define a (*fully*) *parallel* semantics where components of a parallel tuple are always evaluated in parallel regardless of their size. The *mode* of an evaluation is sequential (written seq), parallel (written par), or oracle (written orc). We let  $\alpha$  range over modes:

$$\alpha ::= \text{seq} \mid \text{par} \mid \text{orc}.$$

In addition to evaluating expression, the dynamic semantics also returns cost measures including *raw work* and *raw span* denoted by  $w$  and  $s$  (and variants), and *total work* and *total span*, denoted by  $\mathbb{W}$  and  $\mathbb{S}$  (and variants). Dynamic semantics, whose inductive definition appears in Figure 3, is presented in the style of a natural (big-step) semantics and consists of evaluation judgments of the form

$$e \Downarrow^\alpha v, (w, s), (\mathbb{W}, \mathbb{S}).$$

This judgment states that evaluating expression  $e$  in mode  $\alpha$  yields value  $v$  resulting in raw work of  $w$  and raw span of  $s$  and total work of  $\mathbb{W}$  and total span of  $\mathbb{S}$ .

Figure 3 shows the complete inductive definition of the dynamic cost semantics judgment  $e \Downarrow^\alpha v, (w, s), (\mathbb{W}, \mathbb{S})$ . When evaluating any expression that is not a parallel tuple, we calculate the (raw or total) work and the (raw or total) span by summing up those of the premises (subexpressions) and adding one unit to include the cost of the judgment. For all

$$\begin{array}{c}
\text{(value)} \\
\hline
v \Downarrow^\alpha v, (1, 1), (1, 1) \\
\\
\text{(let)} \\
\frac{e_1 \Downarrow^\alpha v_1, (w_1, s_1), (\mathbb{W}_1, \mathbb{S}_1) \quad e_2[v_1/x] \Downarrow^\alpha v, (w_2, s_2), (\mathbb{W}_2, \mathbb{S}_2)}{(\text{let } x = e_1 \text{ in } e_2) \Downarrow^\alpha v, (w_1 + w_2 + 1, s_1 + s_2 + 1), (\mathbb{W}_1 + \mathbb{W}_2 + 1, \mathbb{S}_1 + \mathbb{S}_2 + 1)} \\
\\
\text{(app)} \\
\frac{(v_1 = \text{fun } f.x.e) \quad e[v_2/x, v_1/f] \Downarrow^\alpha v, (w, s), (\mathbb{W}, \mathbb{S})}{(v_1 \ v_2) \Downarrow^\alpha v, (w + 1, s + 1), (\mathbb{W} + 1, \mathbb{S} + 1)} \\
\\
\begin{array}{cc}
\text{(first)} & \text{(second)} \\
\frac{}{(\text{fst } (v_1, v_2)) \Downarrow^\alpha v_1, (1, 1), (1, 1)} & \frac{}{(\text{snd } (v_1, v_2)) \Downarrow^\alpha v_2, (1, 1), (1, 1)}
\end{array} \\
\\
\text{(case-left)} \\
\frac{e_1[v_1/x_1] \Downarrow^\alpha v, (w, s), (\mathbb{W}, \mathbb{S})}{\text{case } (\text{inl } v_1) \text{ of } \{\text{inl } x_1.e_1, \text{inr } x_2.e_2\} \Downarrow^\alpha v, (w + 1, s + 1), (\mathbb{W} + 1, \mathbb{S} + 1)} \\
\\
\text{(case-right)} \\
\frac{e_2[v_2/x_2] \Downarrow^\alpha v, (w, s), (\mathbb{W}, \mathbb{S})}{\text{case } (\text{inr } v_2) \text{ of } \{\text{inl } x_1.e_1, \text{inr } x_2.e_2\} \Downarrow^\alpha v, (w + 1, s + 1), (\mathbb{W} + 1, \mathbb{S} + 1)} \\
\\
\text{(tuple)} \\
\frac{e_1 \Downarrow^\alpha v_1, (w_1, s_1), (\mathbb{W}_1, \mathbb{S}_1) \quad e_2 \Downarrow^\alpha v_2, (w_2, s_2), (\mathbb{W}_2, \mathbb{S}_2)}{(e_1, e_2) \Downarrow^\alpha (v_1, v_2), (w_1 + w_2 + 1, s_1 + s_2 + 1), (\mathbb{W}_1 + \mathbb{W}_2 + 1, \mathbb{S}_1 + \mathbb{S}_2 + 1)} \\
\\
\text{(ptuple-seq)} \\
\frac{e_1 \Downarrow^{\text{seq}} v_1, (w_1, s_1), (\mathbb{W}_1, \mathbb{S}_1) \quad e_2 \Downarrow^{\text{seq}} v_2, (w_2, s_2), (\mathbb{W}_2, \mathbb{S}_2)}{(|e_1, e_2|) \Downarrow^{\text{seq}} (v_1, v_2), (w_1 + w_2 + 1, s_1 + s_2 + 1), (\mathbb{W}_1 + \mathbb{W}_2 + 1, \mathbb{S}_1 + \mathbb{S}_2 + 1)} \\
\\
\text{(ptuple-par)} \\
\frac{e_1 \Downarrow^{\text{par}} v_1, (w_1, s_1), (\mathbb{W}_1, \mathbb{S}_1) \quad e_2 \Downarrow^{\text{par}} v_2, (w_2, s_2), (\mathbb{W}_2, \mathbb{S}_2)}{(|e_1, e_2|) \Downarrow^{\text{par}} (v_1, v_2), (w_1 + w_2 + 1, \max(s_1, s_2) + 1), (\mathbb{W}_1 + \mathbb{W}_2 + 1 + \tau, \max(\mathbb{S}_1, \mathbb{S}_2) + 1 + \tau)} \\
\\
\text{(ptuple-orc-parallelize)} \\
\frac{w_1 \geq \kappa \wedge w_2 \geq \kappa \quad e_1 \Downarrow^{\text{orc}} v_1, (w_1, s_1), (\mathbb{W}_1, \mathbb{S}_1) \quad e_2 \Downarrow^{\text{orc}} v_2, (w_2, s_2), (\mathbb{W}_2, \mathbb{S}_2)}{(|e_1, e_2|) \Downarrow^{\text{orc}} (v_1, v_2), (w_1 + w_2 + 1, \max(s_1, s_2) + 1), (\mathbb{W}_1 + \mathbb{W}_2 + 1 + \tau + \phi, \max(\mathbb{S}_1, \mathbb{S}_2) + 1 + \tau + \phi)} \\
\\
\text{(ptuple-orc-sequentialize)} \\
\frac{w_1 < \kappa \vee w_2 < \kappa \quad e_1 \Downarrow^{\text{(if } w_1 < \kappa \text{ then seq else orc)}} v_1, (w_1, s_1), (\mathbb{W}_1, \mathbb{S}_1) \quad e_2 \Downarrow^{\text{(if } w_2 < \kappa \text{ then seq else orc)}} v_2, (w_2, s_2), (\mathbb{W}_2, \mathbb{S}_2)}{(|e_1, e_2|) \Downarrow^{\text{orc}} (v_1, v_2), (w_1 + w_2 + 1, s_1 + s_2 + 1), (\mathbb{W}_1 + \mathbb{W}_2 + 1 + \phi, \mathbb{S}_1 + \mathbb{S}_2 + 1 + \phi)}
\end{array}$$

Fig. 3. Dynamic cost semantics

expressions, including parallel tuples, each evaluation step contributes 1 to the raw work or raw span. When calculating total work and total span, we take into account the cost of creating a parallel thread  $\tau$  and the cost of making an oracle decision  $\phi$ .

Evaluation of parallel tuples vary depending on the mode.

- **Sequential mode.** Parallel tuples are treated exactly like sequential tuples: evaluating a parallel tuple simply contributes 1 to the raw and the total work (span), which are computed as the sum of the work (span) of the two branches plus 1. In the sequential mode, raw and total work (span) are the same.
- **Parallel mode.** The evaluation of parallel tuples induces an additional constant cost  $\tau$ . The span is computed as the maximum of the spans of the two branches of the parallel tuple plus 1, and work is computed as the sum of the work of the two branches plus  $\tau$ .
- **Oracle mode.** The scheduling of a parallel tuple depends on the amount of raw work involved in the two branches. If the raw work of each branch is more than  $\kappa$ , then the tuple is evaluated in parallel in the oracle mode. Otherwise, the raw work of at least one branch is less than  $\kappa$ , and the tuple is evaluated sequentially. In this case, the evaluation mode of each branch depends on the work of the branch. If a branch contains more than  $\kappa$  units of raw work, then it is evaluated in oracle mode, otherwise it is evaluated in sequential mode. This switch to sequential mode on small branches ensures that the oracle is not called too frequently during the evaluation of a program.

If the parallel tuple is evaluated sequentially, then the raw/total work and span are both calculated as the sum of the span of the branches plus one. If the parallel tuple is evaluated in parallel, then an extra  $\tau$  is included in the total work and span and the span is computed as the maximum of the span of the two branches.

Theorem 5.1 makes it possible to apply directly the greedy-scheduling theorem to the cost semantics. The basic idea of the proof is to show a correspondence between the cost semantics and dags.

**Theorem 5.1 (Greedy scheduling for the cost semantics)**

Assume  $e \Downarrow^{\text{orc}} v, (w, s), (\mathbb{W}, \mathbb{S})$  to hold for some  $v, w$  and  $s$ . Any greedy scheduler executes the expression  $e$  in no more than  $\frac{\mathbb{W}}{P} + \mathbb{S}$  computations steps on  $P$  processors.

*Proof*

In order to invoke Theorem 4.2, which applies to computation dags, we build the computation dag associated with the execution of the expression  $e$ , including vertices that represent the cost of scheduling. To that end, we describe a recursive algorithm for turning an expression  $e$  with total work  $\mathbb{W}$  and total span  $\mathbb{S}$  into the corresponding computation dag, in which the sum of the weights of the vertices is equal to  $\mathbb{W}$ , and the maximal weight of a path is  $\mathbb{S}$ . The algorithm follows the structure of the derivation that  $e$  has total work  $\mathbb{W}$  and total span  $\mathbb{S}$ .

- If the last rule has zero premises, then  $e$  is an atomic expression and  $\mathbb{W} = \mathbb{S} = 1$ . We build the corresponding dag as a single vertex of unit weight.
- If the last rule has one premise, then  $\mathbb{W}$  takes the form  $\mathbb{W}_1 + 1$  and  $\mathbb{S}$  takes the form  $\mathbb{S}_1 + 1$ . Let  $G_1$  be the dag corresponding to the sub-expression described in the

premise. We build  $G$  by extending  $G_1$  with one unit-weight vertex at the bottom, that is, by sequentially composing  $G_1$  with a dag made of a single vertex.

- Otherwise the last rule has two premises. First, consider the case where  $e$  is a let-expression.  $\mathbb{W}$  takes the form  $\mathbb{W}_1 + \mathbb{W}_2 + 1$  and  $\mathbb{S}$  takes the form  $\mathbb{S}_1 + \mathbb{S}_2 + 1$ . Let  $G_1$  and  $G_2$  be the dags corresponding to the two sub-expressions. We build  $G$  by sequentially composing  $G_1$  with a single unit-weight vertex and then with  $G_2$ .
- Consider now the case of a parallel tuple that is sequentialized.  $\mathbb{W}$  takes the form  $\mathbb{W}_1 + \mathbb{W}_2 + 1 + \phi$  and  $\mathbb{S}$  takes the form  $\mathbb{S}_1 + \mathbb{S}_2 + 1 + \phi$ . Let  $G_1$  and  $G_2$  be the dags corresponding to the two branches. We build  $G$  by sequentially composing a vertex of weight  $1 + \phi$  with the sequential composition of  $G_1$  and  $G_2$ .
- Finally, consider the case of a parallel tuple that is parallelized.  $\mathbb{W}$  takes the form  $\mathbb{W}_1 + \mathbb{W}_2 + 1 + \tau + \phi$  and  $\mathbb{S}$  takes the form  $\max(\mathbb{S}_1, \mathbb{S}_2) + 1 + \tau + \phi$ . Let  $G_1$  and  $G_2$  be the dags corresponding to the two branches. We build  $G$  by sequentially composing a vertex of weight  $1 + \tau + \phi$  with the parallel composition of  $G_1$  and  $G_2$ .

It is straightforward to check that, in each case,  $\mathbb{W}$  and  $\mathbb{S}$  match the sum of the weights of the vertices and the total span of the dag being produced.  $\square$

## 6 Work, Span, and Execution Time Analysis

We analyze the impact of thread-creation overheads on parallel execution time and show how these costs can be reduced dramatically by using our oracle semantics. For our analysis, we first consider an *ideal oracle* that always makes perfectly accurate predictions (about the raw work of expressions) without any overhead (*i.e.*,  $\phi = 0$ ). Such an ideal oracle is unrealistic, because it is practically impossible to determine perfectly accurately the raw work of computations. We therefore consider a realistic oracle that approximates the raw work of computations by performing constant work. Our main result is a theorem that shows that the realistic oracle can reduce the thread-creation overheads to any desired constant fraction of the raw work with some increase in span, which we show to be small for a reasonably broad class of computations.

### 6.1 Ideal Oracle

Theorem 6.1 quantifies the relationships between raw work / raw span and total work / total span for the three possible modes.

**Theorem 6.1 (Work and span)** *Consider an expression  $e$  such that  $e \Downarrow^\alpha v, (w, s), (\mathbb{W}, \mathbb{S})$ . Assume  $\phi = 0$ . The following tight bounds can be obtained for total work and total span, on a machine with  $P$  processors where the cost of creating parallel threads is  $\tau$ .*

$\alpha$	Bound on total work	Bound on total span
seq	$\mathbb{W} = w$	$\mathbb{S} = s = w$
par	$\mathbb{W} \leq (1 + \frac{\tau}{2})w$	$\mathbb{S} \leq (1 + \tau)s$
orc	$\mathbb{W} \leq (1 + \frac{\tau}{\kappa+1})w$	$\mathbb{S} \leq (1 + \max(\tau, \kappa))s$

*Proof*

Results for the sequential mode is straightforward by inspection of the semantics of the source language (Figure 3). The other results can be obtained by specializing the general bounds that we present later in this section (Theorem 6.2 and Theorem 6.3). In what follows, we give examples that attain the bounds for the parallel and the oracle modes.

- For the work in parallel mode, consider an expression consisting only of parallel tuples with  $n$  leaves, and thus  $n - 1$  “internal nodes”. The raw work  $w$  is equal to  $n + (n - 1)$  while the total work  $\mathbb{W}$  is equal to  $n + (n - 1)(1 + \tau)$ . The ratio  $\mathbb{W}/w$  can be rewritten as  $1 + \frac{(n-1)\tau}{2n-1}$ , which tends to  $1 + \frac{\tau}{2}$  as  $n$  grows.

- For the span in parallel mode, we can use the same example. Each parallel tuple accounts for 1 in the raw span but for  $1 + \tau$  in the total span. So, the total span can be as much as  $1 + \tau$  times greater than the raw span.

- For the work in oracle mode, consider an expression with  $n$  nested parallel tuples, where tuples are always nested in the right branch of their parent tuple. The tuples are built on top of expressions that involve  $\kappa$  units of work. In the oracle semantics, all the tuples are executed in parallel. The raw work  $w$  is equal to  $n + (n + 1)\kappa$ , and the total work  $\mathbb{W}$  is equal to  $n(1 + \tau) + (n + 1)\kappa$ . The ratio  $\mathbb{W}/w$  is equal to  $1 + \frac{n\tau}{n(\kappa+1)+\kappa}$ , which tends to  $1 + \frac{\tau}{\kappa+1}$  when  $n$  gets large.

- For the span in oracle mode, in the case  $\tau \geq \kappa$ , we use the same example as for the work. The raw span is  $n + 1$  and the total span is  $n(1 + \tau) + \kappa$ . The ratio  $\mathbb{S}/s$  is equal to  $1 + \frac{n\tau+\kappa-1}{n+1}$ , which approaches  $1 + \tau$  as  $n$  grows.

- For the span in oracle mode, in the case  $\kappa \geq \tau$ , we change the example slightly so that now the tuples are built on leaves that involve just less than  $\kappa$  units of work. In the oracle semantics, all the tuples thus get executed sequentially. In this case the raw span is  $n + \kappa$  and the total span is equal to the total work, which is  $n + (n + 1)\kappa$ . The ratio  $\mathbb{S}/s$  can be expressed as  $1 + \frac{n\kappa}{n+\kappa}$ , which approaches  $1 + \kappa$  as  $n$  grows.  $\square$

This theorem leads to some important conclusions. First, the theorem shows that thread-creation (scheduling) costs matter a great deal. In a parallel evaluation, the total work and total span can be as much as  $\tau$  times larger than the raw span and raw work. This essentially implies that a parallel program can be significantly slower than its sequential counterpart. If  $\tau$  is large compared to the number of processors, then even in the ideal setting, where the number of parallel processors is small relative to  $\tau$ , we may observe no speedups. In fact, it is not uncommon to hear anecdotal evidence of this kind of slowdown in modern computer systems.

Second, the theorem shows that evaluation of a program with an ideal oracle can require as much as  $\frac{\kappa}{2}$  less work than in the parallel mode. This comes at a cost of increasing the span by a factor of up to  $\frac{\kappa}{\tau}$ . Increasing the span of a computation can hurt parallel execution time because many parallel schedulers rely on *parallel slackness*, i.e., the availability of large degree of parallelism to achieve optimal speedups, or  $\frac{w}{p} \gg s$ . Unless done carefully, increasing the span can dramatically reduce parallel slackness. In the common case, however, where there is large amounts of parallel slackness, we can safely increase span by a factor of  $\frac{\kappa}{\tau}$  to reduce the thread-creation overheads. Concretely, if parallel slackness is high, then, in the oracle mode, we can select  $\kappa$  such that parallel slackness is preserved—

total span remains small compared to  $\frac{w}{p}$ , because  $\kappa s$  remains relatively small compared to  $\frac{w}{p}$ —and the total work is reduced approximately by a factor of  $\frac{\kappa}{2}$ .

## 6.2 Realistic Oracles

The analysis that we present above makes two unrealistic assumptions about oracles: 1) that they can accurately predict the raw work for a thread, and 2) that the oracle can make predictions in zero time. Realizing a very accurate oracle in practice is difficult, because it requires determining a priori the execution time of a thread. We therefore generalize the analysis by considering an approximate or realistic oracle that can make errors up to a multiplicative factor  $\mu$  when estimating raw work. For example, an oracle can approximate raw work up to a constant factor of  $\mu = 3$ , *i.e.*, a thread with raw work  $w$  would be estimated to perform raw work between  $\frac{w}{3}$  and  $3w$ . Additionally, we allow the oracle to take some fixed time, written  $\phi$ , to provide its answer.

We show that even with a realistic oracle, we can reduce thread-creation overheads. We start with bounding the span; the result implies that the total span is no larger than  $\mu\kappa$  times the raw span when  $\kappa$  is large compared to  $\tau$  and  $\phi$ .

### Theorem 6.2 (Span with realistic oracle)

$$e \Downarrow^{\text{orc}} v, (w, s), (\mathbb{W}, \mathbb{S}) \quad \Rightarrow \quad \mathbb{S} \leq (1 + \max(\tau, \mu\kappa) + \phi) s$$

#### *Proof*

Let  $\rho$  denote  $1 + \max(\tau, \mu\kappa) + \phi$ ; we want to prove that  $\mathbb{S} \leq \rho s$ . The proof is by induction on the derivation  $e \Downarrow^{\text{orc}} v, (w, s), (\mathbb{W}, \mathbb{S})$ .

- For a rule with zero premise, we have  $\mathbb{S} = s = 1$ . Because  $\rho \geq 1$ , it follows that  $\mathbb{S} \leq \rho s$ .
- For a rule with one premise, we know by induction hypothesis that  $\mathbb{S} \leq \rho s$ . Using again the fact that  $\rho \geq 1$ , we can deduce the inequality  $\mathbb{S} + 1 \leq \rho(s + 1)$ .
- For a rule with two premises that does not correspond to a parallel tuple, we can similarly establish the conclusion  $\mathbb{S}_1 + \mathbb{S}_2 + 1 \leq \rho(s_1 + s_2 + 1)$  using the induction hypotheses  $\mathbb{S}_1 \leq \rho s_1$  and  $\mathbb{S}_2 \leq \rho s_2$ .
- Now, consider the case of a parallel tuple. First, assume that the two branches of this tuple are predicted to be large. In this case, the tuple is executed in parallel and the branches are executed in oracle mode. We exploit the induction hypotheses  $\mathbb{S}_1 \leq \rho s_1$  and  $\mathbb{S}_2 \leq \rho s_2$  to conclude as follows.

$$\begin{aligned} \mathbb{S} &= \max(\mathbb{S}_1, \mathbb{S}_2) + 1 + \tau + \phi \\ &\leq \max(\rho s_1, \rho s_2) + 1 + \max(\tau, \mu\kappa) + \phi \\ &\leq \max(\rho s_1, \rho s_2) + \rho \\ &\leq \rho (\max(s_1, s_2) + 1) \\ &\leq \rho s \end{aligned}$$

- Consider now the case where both branches are predicted to be small. In this case, the tuple is executed sequentially. Because the oracle predicts the branches to be smaller than  $\kappa$ , they must be actually smaller than  $\mu\kappa$ . So, we have  $w_1 \leq \mu\kappa$  and  $w_2 \leq \mu\kappa$ . Moreover, both branches are executed according to the sequential mode, so we have  $\mathbb{S}_1 = w_1$  and  $\mathbb{S}_2 = w_2$ . It follows that  $\mathbb{S}_1 \leq \mu\kappa$  and  $\mathbb{S}_2 \leq \mu\kappa$ . Below, we also exploit the fact that

$\max(s_1, s_2) \geq 1$ , which comes from the fact that raw span is at least one unit. We conclude as follows.

$$\begin{aligned}
\mathbb{S} &= \mathbb{S}_1 + \mathbb{S}_2 + 1 + \phi \\
&\leq \mu\kappa + \mu\kappa + 1 + \phi \\
&\leq (1 + \mu\kappa + \phi) * 2 \\
&\leq (1 + \max(\tau, \mu\kappa) + \phi) \cdot (\max(s_1, s_2) + 1) \\
&\leq \rho s
\end{aligned}$$

• It remains to consider the case where one branch is predicted to be smaller than the cutoff while the other branch is predicted to be larger than the cutoff. In this case again, both branches are executed sequentially. Without loss of generality, assume that the second branch is predicted to be small. In this case, we have  $w_2 \leq \mu\kappa$ . This second branch is thus executed according to the sequential mode, so we have  $\mathbb{S}_2 = s_2 = w_2$ . It follows that  $\mathbb{S}_2 \leq \mu\kappa$ . For the first branch, which is executed according to the oracle mode, we can exploit the induction hypothesis  $\mathbb{S}_1 \leq \rho s_1$ . We conclude as follows:

$$\begin{aligned}
\mathbb{S} &= \mathbb{S}_1 + \mathbb{S}_2 + 1 + \phi \\
&\leq \rho s_1 + \mu\kappa + 1 + \phi \\
&\leq \rho s_1 + (1 + \max(\tau, \mu\kappa) + \phi) \\
&\leq \rho (s_1 + 1) \\
&\leq \rho (\max(s_1, s_2) + 1) \\
&\leq \rho s
\end{aligned}$$

□

This ends our analysis of the span. Now, let us focus on the work. The fact that every call to the oracle can induce a cost  $\phi$  can lead the work to be multiplied by a factor in proportion with  $\phi$ . For example, consider a program made of a complete tree made of  $2n$  leaves,  $n$  parallel tuples on the last layer of the tree, and  $n - 1$  sequential tuples in the upper layers of the tree. This program has raw work equal to  $(n - 1) + n + 2n$ , and total work equal to  $(n - 1) + n\phi + 2n$ . Thus, the increase in work  $\mathbb{W}/w$  is equal to  $\frac{(3+\phi)n-1}{4n-1}$ , which is close to  $\frac{\phi}{4}$  when  $n$  and  $\phi$  are not small. This example shows that a program executed according to the oracle semantics can slow down by as much as a factor  $\frac{\phi}{4}$  compared with a purely sequential execution.

The problem with the above example is that the oracle is called infrequently—only at the leaves of the computation—preventing us from amortizing the cost of the oracle towards larger pieces of computations. Fortunately, most programs do not exhibit this pathological behavior, because parallel tuples are often performed close to the root of the computation, allowing us to detect smaller pieces of work early.

One way to prevent the oracle from being called on smaller pieces of work is to make sure that it is called at regular intervals. For proving a strong bound on the work, we will simply assume that the oracle is not called on small threads by restricting our attention to *balanced programs*. To this end, we define balanced programs as programs that call the oracle only on expressions that are no smaller than some constant  $\gamma$  off from the value  $\frac{\kappa}{\mu}$ , for some  $\gamma \geq 1$ . Note that we use  $\frac{\kappa}{\mu}$  as a target and not  $\kappa$  so as to accommodate possible over-estimations in the estimations of raw work. The formal definition follows.

**Definition 6.1 (Balanced programs)** For  $\gamma \geq 1$ , a program or expression  $e$  is  $\gamma$ -balanced if evaluating  $e$  in the oracle mode invokes the oracle only for subexpressions whose raw work is no less than  $\frac{\kappa}{\mu\gamma}$ .

Note that if a program is  $\gamma$ -balanced and if  $\gamma < \gamma'$ , then this program is also  $\gamma'$ -balanced. We will later give a sufficient condition for proving that particular programs are balanced (Section 6.4). For balanced programs, we are able to bound the total work with respect to the raw work.

**Theorem 6.3 (Work with a realistic oracle)** Let  $e$  be a  $\gamma$ -balanced program.

$$e \Downarrow^{\text{orc}} v, (w, s), (\mathbb{W}, \mathbb{S}) \Rightarrow \mathbb{W} \leq \left(1 + \frac{\mu(\tau + \gamma\phi)}{\kappa + 1}\right) w$$

*Proof*

We establish the following slightly tighter inequality (tighter because  $\gamma \geq 1$  and  $\mu \geq 1$ ).

$$\mathbb{W} \leq \left(1 + \frac{\tau}{\kappa/\mu + 1} + \frac{\phi}{\kappa/(\mu\gamma) + 1}\right) w$$

Define  $\kappa'$  as a shorthand for  $\kappa/\mu$  and  $\kappa''$  as a shorthand for  $\kappa/(\mu\gamma)$ . Note that, because  $\gamma \geq 1$ , we have  $\kappa'' \leq \kappa'$ . Let  $x^+$  be defined as the value  $x$  when  $x$  is nonnegative and as zero otherwise. We establish the following inequality by induction.

$$\mathbb{W} \leq w + \tau \left\lfloor \frac{(w - \kappa')^+}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{(w - \kappa'')^+}{\kappa'' + 1} \right\rfloor$$

This is indeed a strengthened result because we have:

$$\begin{aligned} \tau \left\lfloor \frac{(w - \kappa')^+}{\kappa' + 1} \right\rfloor &\leq \tau \frac{w}{\kappa' + 1} \leq \frac{\tau}{\kappa/\mu + 1} w \\ \text{and } \phi \left\lfloor \frac{(w - \kappa'')^+}{\kappa'' + 1} \right\rfloor &\leq \phi \frac{w}{\kappa'' + 1} \leq \frac{\phi}{\kappa/(\mu\gamma) + 1} w. \end{aligned}$$

The proof is conducted by induction on the derivation of the reduction hypothesis.

- For a rule with zero premises, which describe an atomic operation, we have  $\mathbb{W} = w = 1$ , so the conclusion is satisfied.

- For a rule with a single premise, the induction hypothesis is:

$$\mathbb{W} \leq w + \tau \left\lfloor \frac{(w - \kappa')^+}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{(w - \kappa'')^+}{\kappa'' + 1} \right\rfloor.$$

So, we can easily derive the conclusion:

$$\mathbb{W} + 1 \leq (w + 1) + \tau \left\lfloor \frac{((w + 1) - \kappa')^+}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{((w + 1) - \kappa'')^+}{\kappa'' + 1} \right\rfloor.$$



• For a rule with two premises, we exploit the mathematical inequality  $\left\lfloor \frac{n}{q} \right\rfloor + \left\lfloor \frac{m}{q} \right\rfloor \leq \left\lfloor \frac{n+m}{q} \right\rfloor$ . We have:

$$\begin{aligned} \mathbb{W} &= \mathbb{W}_1 + \mathbb{W}_2 + 1 \\ &\leq w_1 + \tau \left\lfloor \frac{(w_1 - \kappa')^+}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{(w_1 - \kappa'')^+}{\kappa'' + 1} \right\rfloor \\ &\quad + w_2 + \tau \left\lfloor \frac{(w_2 - \kappa')^+}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{(w_2 - \kappa'')^+}{\kappa'' + 1} \right\rfloor + 1 \\ &\leq w + \tau \left\lfloor \frac{(w_1 - \kappa')^+ + (w_2 - \kappa')^+}{\kappa' + 1} \right\rfloor \\ &\quad + \phi \left\lfloor \frac{(w_1 - \kappa'')^+ + (w_2 - \kappa'')^+}{\kappa'' + 1} \right\rfloor. \end{aligned}$$

To conclude, we need to establish the following two mathematical inequalities.

$$\begin{aligned} (w_1 - \kappa')^+ + (w_2 - \kappa')^+ &\leq ((w_1 + w_2 + 1) - \kappa')^+ \\ (w_1 - \kappa'')^+ + (w_2 - \kappa'')^+ &\leq ((w_1 + w_2 + 1) - \kappa'')^+ \end{aligned}$$

The two equalities can be proved in a similar way. Let us establish the first one. There are four cases to consider. First, if both  $w_1$  and  $w_2$  are less than  $\kappa'$ , then the right-hand side is zero, so we are done. Second, if both  $w_1$  and  $w_2$  are greater than  $\kappa'$ , then all the expressions are nonnegative, and we are left to check the inequality  $w_1 - \kappa' + w_2 - \kappa' \leq w_1 + w_2 + 1 - \kappa'$ . Third, if  $w_1$  is greater than  $\kappa'$  and  $w_2$  is smaller than  $\kappa'$ , then the inequality becomes  $(w_1 - \kappa')^+ \leq ((w_1 - \kappa') + (w_2 + 1))^+$ , which is clearly true. The case  $w_1 \geq \kappa'$  and  $w_2 < \kappa'$  is symmetrical. This concludes the proof.

• Consider now the case of a parallel tuple where both branches are predicted to involve more than  $\kappa$  units of work. This implies  $w_1 \geq \kappa'$  and  $w_2 \geq \kappa'$ . In this case, a parallel thread is created. Note that, because  $\kappa'' \leq \kappa'$ , we also have  $w_1 \geq \kappa''$  and  $w_2 \geq \kappa''$ . So, all the values involved in the following computations are nonnegative. Using the induction hypotheses, we have:

$$\begin{aligned} \mathbb{W} &= \mathbb{W}_1 + \mathbb{W}_2 + 1 + \tau + \phi \\ &\leq w_1 + \tau \left\lfloor \frac{w_1 - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w_1 - \kappa''}{\kappa'' + 1} \right\rfloor \\ &\quad + w_2 + \tau \left\lfloor \frac{w_2 - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w_2 - \kappa''}{\kappa'' + 1} \right\rfloor + 1 + \tau + \phi \\ &\leq (w_1 + w_2 + 1) + \tau \left( \left\lfloor \frac{w_1 - \kappa'}{\kappa' + 1} \right\rfloor + \left\lfloor \frac{w_2 - \kappa'}{\kappa' + 1} \right\rfloor + 1 \right) \\ &\quad + \phi \left( \left\lfloor \frac{w_1 - \kappa''}{\kappa'' + 1} \right\rfloor + \left\lfloor \frac{w_2 - \kappa''}{\kappa'' + 1} \right\rfloor + 1 \right) \\ &\leq w + \tau \left\lfloor \frac{(w_1 - \kappa') + (w_2 - \kappa') + (\kappa' + 1)}{\kappa' + 1} \right\rfloor \\ &\quad + \phi \left\lfloor \frac{(w_1 - \kappa'') + (w_2 - \kappa'') + (\kappa'' + 1)}{\kappa'' + 1} \right\rfloor \\ &\leq w + \tau \left\lfloor \frac{(w_1 + w_2 + 1) - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{(w_1 + w_2 + 1) - \kappa''}{\kappa'' + 1} \right\rfloor \\ &\leq w + \tau \left\lfloor \frac{w - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w - \kappa''}{\kappa'' + 1} \right\rfloor. \end{aligned}$$

• Assume now that the two branches are predicted to be less than the cutoff. This implies  $w_1 \leq \kappa'$  and  $w_2 \leq \kappa'$ . Both these threads are executed sequentially, so  $\mathbb{W}_1 = w_1$  and  $\mathbb{W}_2 = w_2$ . Since the program is  $\gamma$ -balanced, we have  $w_1 \geq \kappa''$  and  $w_2 \geq \kappa''$ . Those inequalities ensure that we are able to pay for the cost of calling the oracle, that is, the cost  $\phi$ . Indeed,

since we have  $w_1 + w_2 + 1 - \kappa'' \geq \kappa'' + 1$ , we know that  $\left\lfloor \frac{w_1 + w_2 + 1 - \kappa''}{\kappa'' + 1} \right\rfloor \geq 1$ . Therefore:

$$\begin{aligned} \mathbb{W} &= \mathbb{W}_1 + \mathbb{W}_2 + 1 + \phi \\ &\leq w_1 + w_2 + 1 + \phi \\ &\leq (w_1 + w_2 + 1) + \phi \left\lfloor \frac{w_1 + w_2 + 1 - \kappa''}{\kappa'' + 1} \right\rfloor \\ &\leq w + \tau \left\lfloor \frac{(w - \kappa')^+}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w - \kappa''}{\kappa'' + 1} \right\rfloor. \end{aligned}$$

• It remains to consider the case where one branch is predicted to be bigger than the cutoff while the other is predicted to be smaller than the cutoff. For example, assume  $w_1 \geq \kappa'$  and  $w_2 \leq \kappa'$ . The parallel tuple is thus executed as a sequential tuple. The first thread is executed in oracle mode, whereas the second thread is executed in the sequential mode. For the first thread, we can invoke the induction hypothesis  $\mathbb{W}_1 \leq w_1 + \tau \left\lfloor \frac{w_1 - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w_1 - \kappa''}{\kappa'' + 1} \right\rfloor$ . For the second thread, which is executed sequentially, we have  $\mathbb{W}_2 = w_2$ . Moreover, since the oracle is invoked to predict the size of this second thread, we know by the hypothesis of  $\gamma$ -balance that  $w_2 \geq \kappa / (\mu \gamma) = \kappa''$ . Hence, we have  $\left\lfloor \frac{w_2 + 1}{\kappa'' + 1} \right\rfloor \geq 1$ . We conclude:

$$\begin{aligned} \mathbb{W} &= \mathbb{W}_1 + \mathbb{W}_2 + 1 + \phi \\ &\leq w_1 + \tau \left\lfloor \frac{w_1 - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w_1 - \kappa''}{\kappa'' + 1} \right\rfloor + w_2 + 1 + \phi \\ &\leq w_1 + \tau \left\lfloor \frac{w_1 - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w_1 - \kappa''}{\kappa'' + 1} \right\rfloor + w_2 + 1 + \phi \left\lfloor \frac{w_2 + 1}{\kappa'' + 1} \right\rfloor \\ &\leq w + \tau \left\lfloor \frac{w_1 + w_2 + 1 - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w_1 + w_2 + 1 - \kappa''}{\kappa'' + 1} \right\rfloor \\ &\leq w + \tau \left\lfloor \frac{w - \kappa'}{\kappa' + 1} \right\rfloor + \phi \left\lfloor \frac{w - \kappa''}{\kappa'' + 1} \right\rfloor. \end{aligned}$$

□

We are now ready to combine the version of Brent's theorem adapted to our cost semantics with the bounds that we have established for the total work and span in  $\gamma$ -balanced parallel programs executed under the oracle semantics.

#### Theorem 6.4 (Execution time with a realistic oracle)

Assume an oracle that costs  $\phi$  and makes an error by a factor not exceeding  $\mu$ . Assume  $\kappa > \tau$ , which is always the case in practice. The execution time of a parallel  $\gamma$ -balanced program on a machine with  $P$  processors under the oracle semantics with a greedy scheduler does not exceed the value:

$$\left(1 + \frac{\mu(\tau + \gamma\phi)}{\kappa}\right) \frac{w}{P} + (\kappa\mu + \phi + 1)s.$$

*Proof*

The bound follows by our generalized version of Brent's theorem (Theorem 5.1), and by the bounds established in Theorem 6.3 and Theorem 6.2. For simplicity, we replace the denominator  $\kappa + 1$  with  $\kappa$ . This change does not loosen the bound significantly because  $\kappa$  is usually large compared to the unit cost. Also for simplicity, we have replaced the term  $\max(\tau, \mu\kappa)$  with  $\mu\kappa$ , exploiting the assumption  $\kappa > \tau$  and the fact that  $\mu \geq 1$ . □

### 6.3 Choice of the Cutoff

Theorem 6.4 shows that the running time of a parallel program can be controlled by changing the constant  $\kappa$ . The formula reveals an interesting tradeoff: we can reduce task-creation overheads but this comes at the cost of increasing the span. To see this connection better, consider the bound that appears in the statement of Theorem 6.4 and notice that as the value of  $\kappa$  increases, the work (first) term decreases but the span (second) term increases. The parallel run time thus decreases as we increase  $\kappa$  up to some inflection point and then starts increasing. We compute the optimal value for  $\kappa$  by solving for the root of the derivative. We obtain:

$$\kappa^* = \sqrt{\tau + \gamma\phi} \cdot \sqrt{\frac{w}{Ps}}.$$

Thus, with prior knowledge of the raw work and raw span of a computation, we can pick  $\kappa$  to ensure efficiency of parallel programs.

Such knowledge, however, is often unavailable. As we now show, we can improve efficiency of parallel programs by selecting a fixed  $\kappa$  that guarantees that the task creation overheads can be bounded by any constant fraction of the raw work, without increasing the span of the computation significantly.

**Theorem 6.5 (Run time with fixed  $\kappa$ )** *Consider an oracle with  $\phi$  cost and  $\mu$  error. For any  $\gamma \geq 1$  and for any constant  $r$  such that  $0 < r < 1$ , there exists a constant  $\kappa$  and a constant  $c$  such that the evaluation with the oracle semantics of a  $\gamma$ -balanced program reduces task creation overheads to a fraction  $r$  of the raw work, while in the same time increasing the total span by no more than a factor  $\frac{c}{r}$ . With a greedy scheduler, the total parallel run time on  $P$  processors of such a program thus does not exceed  $(1+r)\frac{w}{P} + \frac{c}{r}s$ .*

*Proof*

Consider a particular  $\gamma$ -balanced program with raw work  $w$  and raw span  $s$ , and consider its evaluation under the oracle semantics. By Theorem 6.3 we know that total work does not exceed  $\left(1 + \frac{\mu(\tau + \gamma\phi)}{\kappa}\right)w$ . To achieve the desired bound on execution time, we take  $\kappa = \frac{\mu(\tau + \gamma\phi)}{r}$ . Plugging this value of  $\kappa$  into the formula yields  $(1+r)w$  for total work, showing that task creation overheads are reduced to a fraction  $r$  of the raw work.

Furthermore, by Theorem 6.2 we know that the total span is bounded by  $(\max(\tau, \mu\kappa) + \phi + 1)s$ . Plugging in the same value for  $\kappa$  yields the following bound on total span:

$$\mathbb{S} \leq \left( \max\left(\tau, \frac{\mu^2(\tau + \gamma\phi)}{r}\right) + \phi + 1 \right) s.$$

Using  $\mu \geq 1$  and  $r < 1$ , we can derive the inequality

$$\mathbb{S} \leq \left( \frac{\mu^2(\tau + \gamma\phi)}{r} + \frac{\phi + 1}{r} \right) s.$$

Choosing  $c = \mu^2(\tau + \gamma\phi) + \phi + 1$  therefore ensures that the total span does not exceed the desired bound  $\frac{c}{r}s$ . The run-time bound follows by an application of our generalized version of Brent's theorem (Theorem 5.1).  $\square$

This final theorem enables us to reduce task creation overheads to any desired constant fraction of the raw work by choosing a  $\kappa$  that is independent of the specific inputs. This

comes at the cost of increasing the span, but only by a constant factor of  $\frac{c}{r}$ . In the common case, when the work is asymptotically greater than span, *e.g.*,  $\Theta(n)$  versus  $O(\log n)$ , the resulting run-time guarantees that the increase in span remain small: specifically, the span term itself is a fraction of the work term for all but a constant number of small inputs.

#### 6.4 Balanced Programs

Our bounds with the realistic oracle hold only for what we called  $\gamma$ -balanced programs, where the oracle is not called on small threads. This assumption can be satisfied by calling the oracle “regularly.” It seems likely that this assumption would hold for many programs without requiring any changes to the program code. In this section, we show that recursive, divide-and-conquer programs are  $\gamma$ -balanced.

To that end, we introduce the notion of  $\gamma$ -regularity. Intuitively, a program is  $\gamma$ -regular if, between any two calls to the oracle involved in the execution of this program, the amount of work does not reduce by more than a factor  $\gamma$ . We will then establish that any  $\gamma$ -regular program is a  $\gamma$ -balanced program. Before giving the formal definition of  $\gamma$ -regularity, we need to formally define what it means for a parallel tuple to be dominated by another parallel tuple.

**Definition 6.2 (Domination of a parallel branch)** *A branch  $e$  of a parallel tuple is said to be dominated by the branch  $e_i$  of another parallel tuple  $(\{e_1, e_2\})$  if the expression  $e$  is involved in the execution of the branch  $e_i$ .*

**Definition 6.3 (Regularity of a parallel program)** *A program is said to be  $\gamma$ -regular if, for any parallel branch involving, say,  $w$  units of raw work, either  $w$  is greater than  $\kappa/(\mu\gamma)$  or this branch is dominated by another parallel branch that involves less than  $\gamma w$  units of work.*

Remark: the condition “ $w$  is greater than  $\kappa/(\mu\gamma)$ ” is typically used to handle the outermost parallel tuples, which are not dominated by any other tuple.

Note that the regularity of a program is always greater than 2. Indeed, if one of the branch of a parallel tuple is more than half of the size of the entire tuple, then the other branch must be smaller than half of that size. On the one hand, algorithms that divide their work in equal parts are  $\gamma$ -regularity with  $\gamma$  very close to 2. On the other hand, ill-balanced programs can have a very high degree of regularity. Observe that every program is  $\infty$ -regular.

For example, consider a program that traverses a complete binary tree in linear time. A call on a tree of size  $n$  has raw work  $nc$ , for some constant  $c$ . If the tree is not a leaf, its size  $n$  has to be at least 3. The next recursive call involves raw work  $\lfloor \frac{n-1}{2} \rfloor c$ . The ratio between those two values is equal  $n / \lfloor \frac{n-1}{2} \rfloor$ . This value is always less than 3 when  $n \geq 3$ . So, the traversal of a complete binary tree is a 3-regular algorithm.

The following lemma explains how the regularity assumption can be exploited to ensure that the oracle is never invoked on threads of size less than  $\kappa/(\mu\gamma)$ . This suggests that, for the purpose of amortizing well the costs of the oracle, a smaller regularity is better.

**Lemma 6.1 (From regularity to balanced)**

*If a program is  $\gamma$ -regular then it is  $\gamma$ -balanced.*

*Proof*

We have to show that, during the execution of a  $\gamma$ -regular program according to oracle semantics, the oracle is never invoked on subexpressions involving less than  $\kappa/(\mu\gamma)$  raw work. Consider a particular subexpression  $e$  involving  $w$  units of raw work, and assume that the oracle is invoked on this subexpression. Because the oracle is being invoked,  $e$  must correspond to the branch of a parallel tuple. By the regularity assumption, either  $w$  is no less than  $\kappa/(\mu\gamma)$ , in which case the conclusion holds immediately, or the branch  $e$  is dominated by a branch  $e_i$  that involves that involves  $w'$  units of work, with  $w' \leq \gamma w$ . For the latter case, we need to establish  $w \geq \kappa/(\mu\gamma)$ . To that end, it suffices to prove that  $w' \geq \kappa/\mu$ , which amounts to showing that the amount of raw work associated with the dominating branch  $e_i$  contains at least  $\kappa/\mu$  raw work.

We conclude the proof by establishing the inequality  $w' \geq \kappa/\mu$ . Because the oracle is being invoked on the subexpression  $e$ , it means that  $e$  is being evaluated in the mode *orc*. Therefore, the call to the oracle on the dominating branch  $e_i$  must have predicted  $e_i$  to contain more than  $\kappa$  raw work. (Otherwise  $e_i$  and its subexpression  $e$  would have both been executed in the sequential mode.) Given that the oracle makes error by no more than a factor  $\mu$ , if  $e_i$  is predicted to contain more than  $\kappa$  units of raw work, then  $e_i$  must contain at least  $\kappa/\mu$  units of raw work. So,  $w' \geq \kappa/\mu$ .  $\square$

## 7 Approximating the Oracle-Guided Semantics

In the previous section, we have established that, if we have access to an oracle that can estimate actual raw work, i.e., sequential run-time, of an expression within a factor of no more than  $\mu$  and at a cost of no more than a constant  $\phi$ , then we can effectively control the thread-creation overheads and thus enable efficient parallel execution. In this section, we describe how to approximate such an oracle by combining certain information from the programmer along with run-time measurements. We first describe our approximation algorithm and then, in Section 8, show how some of the information needed by the algorithm can be derived via a program transformation.

The basic idea behind our approximation algorithm is to use some crucial information provided by the programmer to estimate the actual raw work of a computation. More precisely, we require the user to provide a *cost function* for each function in the program. When applied to an argument  $v$ , a cost function of  $f$  returns an *abstract cost* of the raw-work of application of  $f$  to  $v$ . A crucial property of the abstract costs is that they should be abstract enough that the programmer can write the cost functions without necessarily knowing the details of the hardware where the programs will be executed on. Yet, abstract costs should provide sufficient information to estimate the actual run times.

Asymptotic complexity specifications serve as a natural cost function by satisfying both of these properties. Since they eliminate hardware specific constants, they can be specified easily. Using complexity functions, we can approximate the actual run time of sequentially executed functions by simply determining the constants hidden by the asymptotic complexity notation. Such an approximation can be performed, e.g., by using the least squares method or similar techniques for data fitting from known samples. To perform this approximation we use an *estimator* data structure, that, given abstract cost samples for a function, can estimate the raw work for that function on a given argument.

```
signature Estimator =  
sig  
  type cost  
  type estimator  
  val create: unit → estimator  
  val report: estimator × cost × float → unit  
  val predict: estimator × cost → float  
end
```

Fig. 4. The interface for the estimator data structure.

Figure 4 shows the interface for the estimator. To perform accurate estimates, the estimator utilizes profiling data obtained from actual execution times. The sampling operation `report(t, c, d)` adds a cost `c` and an execution time `d` to the set of samples in an estimator `t`. An estimate of the actual execution time is obtained by calling `predict`. Given an estimator `t` and cost `c`, the call `predict(t, c)` returns a predicted execution time.

Since the abstract cost simply a measure of the asymptotic work of a function, all that remains for the estimator data structure is to calculate at run-time the actual constant factors for the hardware, where execution takes place. In our implementation, we represent a value of type `cost` as an integer that represents the application of the complexity function applied to the input size. We approximate the actual run time by calculating a single constant, assuming that the constants in all terms of the asymptotic complexity are the same. Although assuming a single constant can decrease the precision of the approximations, we believe that it suffices because we only have to compute lower bounds for our functions; i.e., we only need to determine whether they are “big enough” for parallel execution. We note that, for our theoretical bounds to apply, complexity expressions should require constant time to evaluate, which is usually the case because, in the common case, the cost functions are relatively simple functions.

To ensure lightweight execution, our implementation simply computes the constant factor for each reported cost, and averages such constants over a period of time, which it then uses for prediction. The constants calculated for each reported cost may evolve over time. For example, if the current program is sharing the machine with another program, a series of memory reads by the other program may slow down the current program. For this reason, we do not just compute the average across the entire history, but instead maintain a moving average, that is, an average of the values gathered across a certain number of runs.

Computing averages is not entirely straightforward. On the one hand, storing the average in a memory cell that is shared by all processors is not satisfying, because it would involve some synchronization problems. On the other hand, using a different memory cell for every processor is not satisfying either, because it leads to slower updates of the constants when they change. In particular, in the beginning of the execution of a program it is important that all processors quickly share a relatively good estimate of the constant factors.

For these reasons, we have opted for an approach that uses not only a shared memory cell but also one data structure local to every processor. The shared memory cell associated with each estimator contains the estimated value for the constant that is read by all the processors when they need to predict execution times. The local data structures are used to accumulate statistics on the value of the constant. Those statistics are reported on a regular basis to the shared memory cell, by computing a weighted mean between the value previously

stored in the shared memory cell and the value obtained out of the local data structure. We treat initializations somewhat specially: for the first few measures, a processor always begins by reporting its current average to the shared memory cell. This policy ensures a fast propagation of the information gathered from the first runs, so as to quickly improve the accuracy of the predictions.

In the rest of this section we present a more detailed description of our implementation of estimators. An estimator is represented as (1) a shared floating-point value storing the global estimated value, (2) an integer value storing, in the early phases, the number of times that the shared value has been updated, and (3) a processor-indexed array of pairs, each made of an integer and a floating-point value, for storing the number of measures and the sum of the constant factors measured. We define the *local average value* associated with a given processor as the ratio between the sum of the measured constants and the number of those measures.

The function `create` allocates a fresh estimator. It sets all fields to zero, except the global constant which is set to a very pessimistic constant, e.g., 1 microsecond. As a result, in the first few calls, the estimator largely over-approximates the work, leading only very small pieces of computation to be executed sequentially.

The function `predict` simply computes the product of the global estimated value of the constant with the integer cost produced by the user-provided asymptotic cost function.

The function `report` refines the estimation of the constant factor. First, it updates the pair of values associated with the processor that performed the measure, by incrementing the number of measures and adding to the sum field the ratio between the execution time and the complexity value being reported. Second, it decides whether the local average value of the constant should be reported to the global value. It does so if either (a) the number of measures exceeds  $c * P$ , for some real constant  $c$ , or if (b) the number of times that the shared value has been updated is less than a fixed bound. The latter condition helps for fast propagation of the constant in the early phases.

When the function `report` decides to report its local average value to the global value, it writes into the shared value field a weighted average of the local average value and of the previous global value. For safety (that is, to avoid potentially sequentializing too large pieces of computations), we limit a single change to the global value to be of at most one order of magnitude (e.g., a factor 10). Besides, as an optimization, we skip the write operation if we notice that the new value is not significantly different from the old (e.g., when the value would change by less than 20%). After updating the shared value, the function resets to zero the pair of values associated with the current processor, in order to begin a new round of measurements. Moreover, in case the update of the shared value has been triggered by the condition (b) described above, the function also increments the shared field that stores the number of updates to the global value.

Note that the write operation that updates the shared value may sometimes get discarded as a result of a data race between processors. Such races, however, are benign.

When implementing the oracle, we faced three technical difficulties. First, we had to pay attention to the fact that the memory cells allocated for the different processors are not allocated next to each other. Otherwise, those cells would fall in the same cache line, in which case writing in one of these cells would make the other cells be removed from caches, making subsequent reads more costly. Second, we observed that the time

measures typically yield a few outliers. Those are typically due to the activity of the garbage collector or of another program being scheduled by the operating system on the same processor. Fortunately, we have found detecting these outliers to be relatively easy because the measured times are at least one or two orders of magnitude greater than the cutoff value. Third, the default system function that reports the time is only accurate by one microsecond. This is good enough when the cutoff is greater than 10 microseconds. However, if one were to aim for a smaller cutoff, which could be useful for programs exhibiting only a limited amount of parallelism, then more accurate techniques would be required, for example using the specific processor instructions for counting the number of processor cycles.

## 8 Implementation via Source-to-Source Translation

We describe how to compile parallel codes with complexity annotations into codes that implement our oracle-guided semantics. Our compilation technique performs a source-to-source translation and relies on the technique described in Section 7. The idea is to compile every piece of code in two versions: one for the sequential semantics (seq mode), where all parallel pairs are simply erased to sequential pairs, and one for the oracle-guided semantics (orc mode), where parallel pairs are instrumented in such a way as to perform predictions and possibly measures, and to decide whether to execute the branches sequentially or in parallel.

For simplicity, we assume that constituents of parallel tuples are function applications, *i.e.*, they are of the form  $(|f_1 v_1, f_2 v_2|)$ . Note that this assumption does not cause loss of expressiveness, because a term  $e$  can always be replaced by a trivial application of a “think”, a function that ignores its argument (typically of type “unit”) and evaluates  $e$  to a dummy argument. Throughout, we write “`fun f.x.eb [ec]`” to denote a function “`fun f.x.eb`” for which the cost function for the body  $e_b$  is described by the expression  $e_c$ . This expression  $e_c$ , which may refer to the argument  $x$ , should be an expression whose evaluation always terminates and produces a cost of type `cost`.

To associate an estimator with each function, in a simple pass over the source code, we allocate and initialize an estimator for each syntactic function definition. For example, if the source code contains a function of the form “`fun f.x.eb [ec]`”, then our compiler allocates an estimator specific to that function definition. Specifically, if the variable  $r$  refers to the allocated estimator, then the translated function, written “`fun f.x.eb [ec|r]`”, is annotated with  $r$ .

The second pass of our compilation scheme uses the allocated estimators to approximate the actual raw work of function applications and relies on an `MakeBranch` function to determine whether an application should be run in the oracle or in the sequential mode. Figure 5 defines more precisely the second pass. We write  $\llbracket v \rrbracket$  for the translation of a value  $v$ , and we write  $\llbracket e \rrbracket^\alpha$  for the translation of the expression  $e$  according to the semantics  $\alpha$ , which can be either `seq` or `orc`. When specifying the translation, we use triples, quadruples, projections, sequence, if-then-else statements, and unit value; these constructions can all be easily defined in our core programming language.

Translation of values other than functions does not depend on the mode and is relatively straightforward. We translate functions, which are of the form “`fun f.x.eb [ec|r]`”, into a



$\llbracket x \rrbracket$	$\equiv x$
$\llbracket (v_1, v_2) \rrbracket$	$\equiv (\llbracket v_1 \rrbracket, \llbracket v_2 \rrbracket)$
$\llbracket \text{inl } v \rrbracket$	$\equiv \text{inl } \llbracket v \rrbracket$
$\llbracket \text{inr } v \rrbracket$	$\equiv \text{inr } \llbracket v \rrbracket$
$\llbracket \text{fun } f.x.e_b [e_c   r] \rrbracket$	$\equiv (r, (\text{fun } \_x. \llbracket e_c \rrbracket^{\text{seq}}), (\text{fun } f.x. \llbracket e_b \rrbracket^{\text{seq}}), (\text{fun } f.x. \llbracket e_b \rrbracket^{\text{orc}}))$
$\llbracket v \rrbracket^\alpha$	$\equiv \llbracket v \rrbracket$
$\llbracket v_1 v_2 \rrbracket^{\text{seq}}$	$\equiv \text{proj}^3 \llbracket v_1 \rrbracket \llbracket v_2 \rrbracket$
$\llbracket v_1 v_2 \rrbracket^{\text{orc}}$	$\equiv \text{proj}^4 \llbracket v_1 \rrbracket \llbracket v_2 \rrbracket$
$\llbracket (e_1, e_2) \rrbracket^\alpha$	$\equiv (\llbracket e_1 \rrbracket^\alpha, \llbracket e_2 \rrbracket^\alpha)$
$\llbracket \text{let } x = e_1 \text{ in } e_2 \rrbracket^\alpha$	$\equiv \text{let } x = \llbracket e_1 \rrbracket^\alpha \text{ in } \llbracket e_2 \rrbracket^\alpha$
$\llbracket \text{fst } v \rrbracket^\alpha$	$\equiv \text{fst } \llbracket v \rrbracket$
$\llbracket \text{snd } v \rrbracket^\alpha$	$\equiv \text{snd } \llbracket v \rrbracket$
$\llbracket \text{case } v \text{ of } \{ \text{inl } x.e_1, \text{inr } x.e_2 \} \rrbracket^\alpha$	$\equiv \text{case } \llbracket v \rrbracket \text{ of } \{ \text{inl } x. \llbracket e_1 \rrbracket^\alpha, \text{inr } x. \llbracket e_2 \rrbracket^\alpha \}$
$\llbracket ( f_1 v_1, f_2 v_2 ) \rrbracket^{\text{seq}}$	$\equiv (\text{proj}^3 \llbracket f_1 \rrbracket \llbracket v_1 \rrbracket, \text{proj}^3 \llbracket f_2 \rrbracket \llbracket v_2 \rrbracket)$
$\llbracket ( f_1 v_1, f_2 v_2 ) \rrbracket^{\text{orc}}$	$\equiv \begin{cases} \text{let } (b_1, k_1) = \text{MakeBranch}(\llbracket f_1 \rrbracket, \llbracket v_1 \rrbracket) \text{ in} \\ \text{let } (b_2, k_2) = \text{MakeBranch}(\llbracket f_2 \rrbracket, \llbracket v_2 \rrbracket) \text{ in} \\ \text{if } (b_1 \ \&\& \ b_2) \text{ then } ( k_1 \ (), k_2 \ () ) \ \text{else } (k_1 \ (), k_2 \ ()) \end{cases}$

Fig. 5. Translation for oracle scheduling

quadruple consisting of the estimator  $r$ , a sequential cost function, the sequential version of the function, and the oracle versions of the function. Translation of a function application depends on the mode. In the sequential mode, the sequential version of the function is selected (by projecting the third component of the function) and used in the application. Similarly, in the oracle mode, the oracle version of the function is selected and used in the application. To translate a tuple, we recursively translate the subexpression, while preserving the mode. Similarly, translation of the `let`, projections, and `case` constructs are entirely structural.

In the sequential mode, a parallel tuple is turned into a simple tuple. In the oracle mode, the translation applies the oracle-based scheduling policy with the aid of the meta-function `MakeBranch`. This meta-function, shown in Figure 6, describes the template of the code generated for preparing the execution of a parallel tuple. `MakeBranch` expects a (translated) function  $f$  and its (translated) argument  $v$ , and it returns a boolean  $b$  indicating whether the application of  $f$  to  $v$  is expected to take more or less time than the cutoff  $\kappa$ , and a thunk  $t$  to execute this application. On the one hand, if the application is predicted to take more time than the cutoff (in which case  $b$  is true), then the thunk  $t$  corresponds to the application of the oracle-semantics version of the function  $f$ . On the other hand, if the application is predicted to take less time than the cutoff (in which case  $b$  is false), then the thunk  $t$  corresponds to the application of the sequential-semantics version of the function  $f$ . Moreover, in the latter case, the time taken to execute the application sequentially is measured. This time measure is reported to the estimator by the auxiliary meta-function `MeasuredRun` (Figure 6), so as to enable its approximations.

Observe that the translation introduces many quadruples and applications of projection functions. However, in practice, the quadruples typically get inlined so most of the projections can be computed at compile time. Observe also that the compilation scheme involves

```

MakeBranch( $f, v$ )  $\equiv$ 
  let  $r = \text{proj}^1 f$  in
  let  $m = \text{proj}^2 f v$  in
  let  $b = \text{predict}(r, m) > \kappa$  in
  let fun  $k_{\text{seq}} () = \text{proj}^3 f v$  in
  let fun  $k'_{\text{seq}} () = \text{MeasuredRun}(r, m, k_{\text{seq}})$  in
  let fun  $k_{\text{orc}} () = \text{proj}^4 f v$  in
  let  $k = \text{if } b \text{ then } k_{\text{orc}} \text{ else } k'_{\text{seq}}$  in
  ( $b, k$ )

MeasuredRun( $r, m, k$ )  $\equiv$ 
  let  $t = \text{get\_time} ()$  in
  let  $v = k ()$  in
  let  $t' = \text{get\_time} ()$  in
  report( $r, m, (t' - t)$ );
  v

```

Fig. 6. Auxiliary meta-functions used for compilation

some code duplication, because every function is translated once for the sequential mode and once for the oracle mode. In theory, the code could grow exponentially when the code involves functions defined inside the body of other functions. In practice, the code the growth is limited because functions are rarely deeply nested. If code duplication was a problem, then we can use flattening to eliminate deep nesting of local functions, or pass the mode  $\alpha$  as an extra argument to functions.

## 9 Empirical Evaluation

In this section, we evaluate the effectiveness of our implementation through several experiments. We consider results from a range of benchmarks run on a machine with sixteen processors. The results show that, in each case, our oracle implementation improves on the plain work-stealing implementation. Furthermore, the results show that the oracle implementation scales well with up to sixteen processors.

### 9.1 Implementation in Manticore

In this section, we describe the implementation of our scheduling technique in an actual language and system. In our approach, source programs are written in our own dialect of the Caml language (37), which is a strict functional language. Our Caml dialect corresponds to the core Caml language extended with syntax for parallel pairs and complexity annotations.

We use the Caml type checker to obtain a typed syntax tree, on which we perform the oracle-scheduling translation defined in Figure 5. We then produce code in the syntax of Parallel ML (PML) (28), a parallel language close to Standard ML. The translation from Caml to PML is straightforward because the two languages are relatively similar. We compile our source programs to x86-64 binaries using Manticore, which is the optimizing PML compiler.

The original theorem of Brent as well as our generalization both assume a greedy scheduler that can find available work (threads to execute) immediately with no overhead. This

assumption is unrealistic of course in a literal sense, but work-stealing schedulers are able to control the overheads for balancing work remarkably well, both in theory (19; 10; 5) and in practice (30). On the one hand, overheads related to queuing threads are tiny, and, on the other hand, the delay due to steals is tamed by the fact that steals are relatively rare. Work stealing has also been shown to have good locality behavior, which further improves its effectiveness (3).

In Manticore's work-stealing scheduler, all system processors are assigned to collaborate on the computation. Each processor owns a deque (doubly-ended queue) of threads represented as thunks. Processors treat their own deques like call stacks. When a processor starts to evaluate a parallel-pair expression, it creates a thread for the second subexpression of the pair and pushes the thread onto the bottom of the deque. Processors that have no work left try to *steal* threads from others. More precisely, they repeatedly select a random processor and try to pop a thread from this processor's deque.

Manticore's implementation of work stealing (45) adopts a code-specialization scheme, called clone translation, taken from Cilk-5's implementation (30).<sup>2</sup> With clone translation, each parallel-pair expression is compiled into two versions: the fast clone and the slow clone. The purpose of a fast clone is to optimize the code that corresponds to evaluating on the local processor, whereas the slow clone is used when the second branch of a parallel-pair is migrated to another processor. A common aspect between clone translation and our oracle translation (Figure 5) is that both generate specialized code for the sequential case. But the clone translation differs in that there is no point at which parallelism is cut off entirely, as the fast clone may spawn threads.

The scheduling cost involved in the fast clone is a (small) constant, because it involves just a few local operations, but the scheduling cost of the slow clone is variable, because it involves inter-processor communication. It is well established, both through analysis and experimentation, that (with high probability) no more than  $O(P\mathbb{S})$  steals occur during the evaluation (10). So, for programs that exhibit parallel slackness ( $\mathbb{W} \gg P\mathbb{S}$ ), we do not need to take into account the cost of slow clones because there are relatively few of them. We focus only on the cost of creating fast clones, which correspond to the cost  $\tau$ . A fast clone needs to package a thread, push it onto the deque and later pop it from the deque. So, a fast clone is not quite as fast as the corresponding sequential code. The exact slowdown depend on the implementation, but in our case we have observed that a fast clone is 3 to 5 times slower than a simple function call.

## 9.2 Test Machine

Our test machine has four quad-core AMD Opteron 8380 processors running at 2.5GHz. Each core has 64Kb each of L1 instruction and data cache, and a 512Kb L2 cache. Each processor has a 6Mb L3 cache that is shared with the four cores of the processor. The system has 32Gb of RAM and runs Debian Linux (kernel version 2.6.31.6-amd64).

<sup>2</sup> In the Cilk-5 implementation, it is called clone compilation.

### 9.3 Measuring Scheduling Costs

We report estimates of the thread-creation overheads for each of our test machine. To estimate, we use a synthetic benchmark expression  $e$  whose evaluation sums integers between zero and 30 million using a parallel divide-and-conquer computation. We chose this particular expression because most of its evaluation time is spent evaluating parallel pairs.

First, we measure  $w_s$ : the time required for executing a sequentialized version of the program (a copy of the program where parallel tuples are systematically replaced with sequential tuples). This measure serves as the baseline. Second, we measure  $w_w$ : the time required for executing the program using work stealing, on a single processor. This measure is used to evaluate  $\tau$ . Third, we measure  $w_o$ : the time required for executing a version of the program with parallel tuples replaced with ordinary tuples but where we still call the oracle. This measure is used to evaluate  $\phi$ .

We then define the work-stealing overhead  $c_w = \frac{w_w}{w_s}$ . We estimate the cost  $\tau$  of creating a parallel thread in work stealing by computing  $\frac{w_w - w_s}{n}$ , where  $n$  is the number of parallel pairs evaluated in the program. We also estimate the cost  $\phi$  of invoking the oracle by computing  $\frac{w_o - w_s}{m}$ , where  $m$  is the number of times the oracle is invoked.

Our measures are as follows:  $c_w = 4.86$  and  $\tau = 0.09$  microseconds and  $\phi = 0.18$  microseconds. The first value indicates that work stealing alone can induce a slowdown by a factor of 4 or 5, for programs that create a huge number of parallel tuples. The value of  $\tau$ , close to one tenth of a microsecond, indicates, that the cost of creating parallel threads is roughly between 200 and 300 processor cycles (since our benchmark machine runs at 2.5GHz). This cost is quite significant in front of basic operations that execute in just a few cycles.

The oracle cost  $\phi$  is even larger than  $\tau$ —about twice larger. The fact that the oracle is associated with significant costs is not surprising: recall that the estimations and measures associated with our implementation of the oracle involve functions calls and memory operations. However, also keep in mind that our investment in performing oracle predictions is supposed to pay off as soon as we are able to sequentialize large pieces of computation. Indeed, as we have established, the cost of the oracle can be expected to be well amortized and to only account for a tiny fraction of the total execution time.

In our experiments, we used the cutoff value  $\kappa = 26$  microseconds. We use Theorem 6.4 to estimate an upper-bound on the overheads associated with thread creation and oracle predictions. To that end, we first need to evaluate  $\mu$  and  $\gamma$ . We can estimate  $\mu$ , the bound on the relative error associated with the predictions, by performing runs where we logged both our predictions and our measures. We observed that, apart from a few outliers (i.e. exceptionally large measures for sequentialized threads), our oracle is always accurate within a factor 2. So, we assume  $\mu = 2$ . Besides, our benchmark programs are fairly regular: we can assume  $\gamma = 3$  for all of them.

Now, Theorem 6.4 tells us that, for programs exhibiting sufficient parallelism, the overheads are essentially bounded by  $\frac{\mu(\tau + \gamma\phi)}{\kappa}$ . In our setting, this ratio evaluates to just below 5%—we actually set the value of  $\kappa$  to meet this target. In practice, since the errors associated with our time predictions tend to balance out, the  $\mu$  term that appears in front of the bound does not impact us as much as it could in theory. Likewise, the irregularity

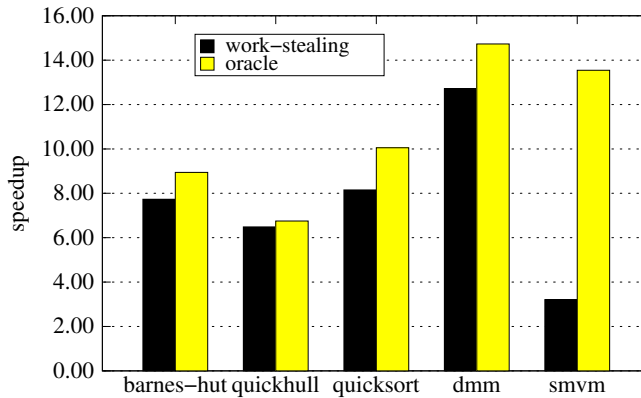


Fig. 7. Comparison of the speedup on sixteen processors. Higher bars are better.

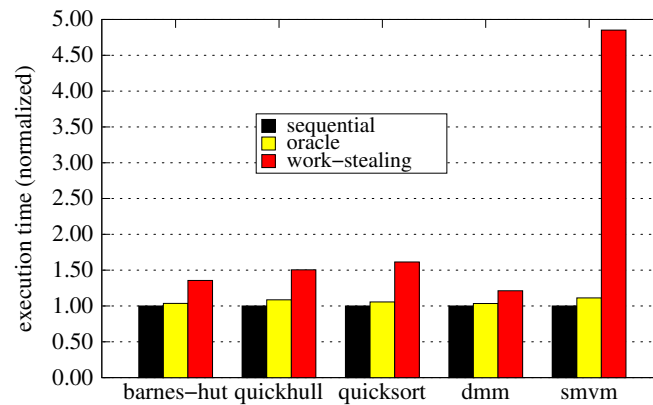


Fig. 8. Comparison of execution times (normalized) on a single processor. Lower bars are better.

factor  $\phi$  does not actually reach 3 for every parallel tuple executed. For these reasons, we believe that, in practice, overheads usually do not exceed 3%.

#### 9.4 Benchmarks

We used five benchmarks in our empirical evaluation. Each benchmark program was originally written by other researchers and ported to our dialect of Caml.

The Quicksort benchmark sorts a sequence of 2 million integers. Our program is adapted from a functional, tree-based algorithm (13). The algorithm runs with  $O(n \log n)$  raw work and  $O(\log^2 n)$  raw depth, where  $n$  is the length of the sequence. Sequences of integers are represented as binary trees in which sequence elements are stored at leaf nodes and each internal node caches the number of leaves contained in its subtree.

The Quickhull benchmark calculates the convex hull of a sequence of 3 million points contained in 2-d space. The algorithm runs with  $O(n \log n)$  raw work and  $O(\log^2 n)$  raw depth, where  $n$  is the length of the sequence. The representation of points is similar to that of Quicksort, except that leaves store 2-d points instead of integers.

The Barnes-Hut benchmark is an  $n$ -body simulation that calculates the gravitational forces between  $n$  particles as they move through 2-d space (11). The Barnes-Hut computation consists of two phases. In the first, the simulation volume is divided into square cells via a quadtree, so that only particles from nearby cells need to be handled individually and particles from distant cells can be grouped together and treated as large particles. The second phase calculates gravitational forces using the quadtree to accelerate the computation. The algorithm runs with  $O(n \log n)$  raw work and  $O(\log n)$  raw depth. Our benchmark runs 10 iterations over 100,000 particles generated from a random Plummer distribution (44). The program is adapted from a Data-Parallel Haskell program (42). The representation we use for sequences of particles is similar to that of Quicksort.

The SMVM benchmark multiplies an  $m \times n$  matrix with an  $n \times 1$  dense vector. Our sparse matrix is stored in the compressed sparse-row format. The program contains parallelism both between dot products and within individual dot products. We use a sparse matrix of dimension  $m = 500,000$  and  $n = 448,000$ , containing 50,400,000 nonzero values.

The DMM benchmark multiplies two dense, square  $n \times n$  matrices using the recursive divide-and-conquer algorithm of Frens and Wise (29). We have recursion go down to scalar elements. The algorithm runs with  $O(n^3)$  raw work and  $O(\log n)$  raw depth. We selected  $n = 512$ .

### 9.5 Performance

For every benchmark, we measure several values.  $T_{\text{seq}}$  denotes the time to execute the sequential version of the program. We obtain the sequential version of the program by replacing each parallel tuple with an ordinary tuple and erasing complexity functions, so that the sequential version includes none of the thread-creation overheads.  $T_{\text{par}}^P$  denotes the execution time with work stealing on  $P$  processors.  $T_{\text{orc}}^P$  denotes the execution time of our oracle-based work stealing on  $P$  processors.

The most important results of our experiments come from comparing plain work stealing and our oracle-based work stealing side by side. Figure 7 shows the speedup on sixteen processors for each of our benchmarks, that is, the values  $T_{\text{par}}^{16}/T_{\text{seq}}$  and  $T_{\text{orc}}^{16}/T_{\text{seq}}$ . The speedups show that, on sixteen cores, our oracle implementation is always between 4% and 76% faster than work stealing.

The fact that some benchmarks benefit more from our oracle implementation than others is explained by Figure 8. This plot shows execution time for one processor, normalized with respect to the sequential execution times. In other words, the values plotted are 1,  $T_{\text{orc}}^1/T_{\text{seq}}$  and  $T_{\text{par}}^1/T_{\text{seq}}$ . The values  $T_{\text{orc}}^1/T_{\text{seq}}$  range from 1.03 to 1.13 (with an average of 1.07), indicating that the thread-creation overheads in the oracle implementation do not exceed 13% of the raw work in any benchmark. The cases where we observe large improvements in speedup are the same cases where there is a large difference between sequential execution time and plain work-stealing execution time. When the difference is large, there is much room for our implementation to improve on work stealing, whereas when the difference is small we can only improve the execution time by a limited factor.

Figure 9 shows speedup curves for each of our experiments, that is, values of  $T_{\text{par}}^P/T_{\text{seq}}$  and  $T_{\text{orc}}^P/T_{\text{seq}}$  against the number of processors  $P$ . The curves show that our oracle implementation generally scales well up to sixteen processors. There is one exception, which

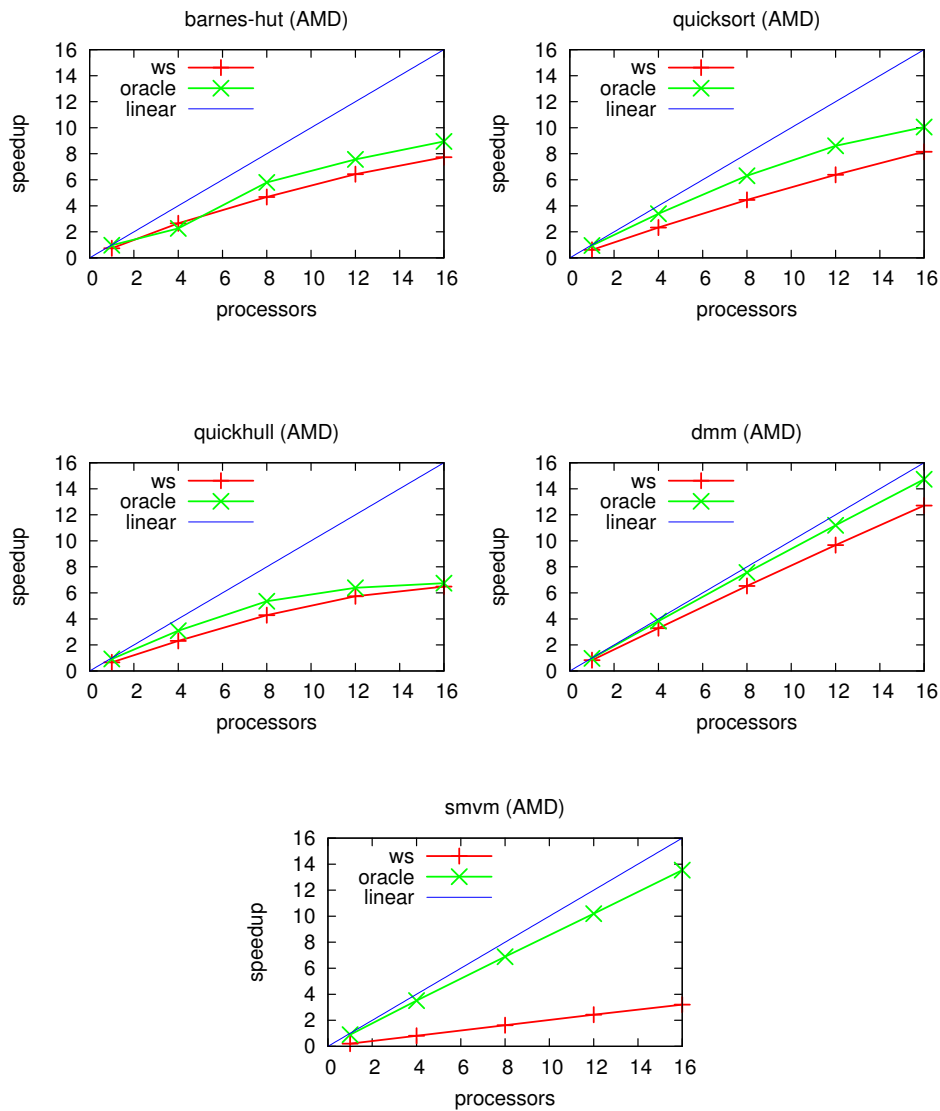


Fig. 9. Comparison between work stealing and oracle.

is the quickhull benchmark. For this benchmark, the curve tails off after reaching twelve processors. We believe that this tailing is either due to a lack of parallelism in the program or to a bottleneck on memory accesses. Notice, however, that our scheduler does not fall below plain work stealing.

## 10 Related Work

### 10.1 Cutting Off Excess Parallelism

This study is not the first to propose using cost prediction to determine when to cut off parallelism. One approach, developed in early work in functional programming, uses list size to determine cut offs (35). Using list size alone is limited, because the technique assumes linear work complexity for every parallel operation.

Another way to handle cost prediction is to use the depth and height of the recursion tree (54; 41). But depth and height are not, in general, the most direct means to predict the execution time of subcomputations. In our oracle scheduling, we ask for either the programmer or compiler to provide for each function a cost function that expresses the asymptotic cost of applying the function.

Lopez *et al.* take this approach as well, but in the context of logic programming (38). On the surface, their technique is similar to our oracle scheduling, except that their cost estimators do not utilize profiling to estimate constant factors. An approach without constant-factor estimation is overly simplistic for modern processors, because it relies on complexity function predicting execution time exactly. On modern processors, execution time depends heavily on factors such as caching, pipelining, *etc.* and it is not feasible in general to predict execution time from a complexity function alone.

### 10.2 Reducing Thread-Creation Cost

One approach to the granularity problem is to focus on reducing the cost of creating a thread, rather than limiting the number of threads created. This approach is taken by implementations of work stealing with lazy thread creation (39; 25; 26; 30; 45; 34; 47; 52). In lazy thread creation, the work stealing scheduler is implemented so as to avoid, in the common case, the major scheduling costs, in particular, those of inter-processor communication. But, in even the most efficient lazy thread creation, there is still a non-negligible scheduling cost for each implicit thread.

Lazy Binary Splitting (LBS) is an improvement to lazy thread creation that applies to parallel loops (52). The crucial optimization comes from extending the representation of a thread so that multiple loop iterations can be packed into a single thread. This representation enables the scheduler to both avoid creating closures and executing deque operations for most iterations. A limitation of LBS is that it addresses only parallel loops whose iteration space is over integers. Lazy Tree Splitting (LTS) generalizes LBS to handle parallel aggregate operations that produce and consume trees, such as map and reduce (12). LTS is limited, however, by the fact that it requires a special cursor data structure to be defined for each tree data structure.

### 10.3 Amortizing Per-Thread Costs

Aharoni *et al.* study the granularity problem in the setting of distributed computing (8), where the crucial issue is how to minimize the cost of inter-processor communication. In their setting, the granularity problem is modeled as a staging problem, in which there are two stages. The first stage consists of a set of processor-local thread pools and the second



stage consists of a global thread pool. Moving a thread to the global thread pool requires inter-processor communication. The crucial decision is how often each processor should promote threads from its local thread pool to the global thread pool. We consider a different model of staging in which there is one stage for parallel evaluation and one for sequential evaluation.

The approach proposed by Aharoni *et al.* is based on an online algorithm called CG. In this approach, it is assumed that the cost of moving a thread to the global thread pool is an integer constant, called  $g$ . The basic idea is to use amortization to reduce the scheduling total cost of moving threads to the global thread pool. In particular, for each thread that is moved to the global thread pool, CG ensures that there are at least  $g + 1$  threads added to the local thread pool. Narlikar describes a similar approach based on an algorithm called DFDeques (40). Just as with work stealing, even though the scheduler can avoid the communication costs in the common case, the scheduler still has to pay a non-negligible cost for each implicit thread.

#### 10.4 Flattening and Fusion

Flattening is a well-known program transformation for nested parallel languages (18). Implementations of flattening include NESL (17) and Data Parallel Haskell (42). Flattening transforms the program into a form that maps well onto SIMD architectures. Flattened programs are typically much simpler to schedule at run time than nested programs, because much of the schedule is predetermined by the flattening (50). Controlling the granularity of such programs is correspondingly much simpler than in general. A limitation of existing flattening is that certain classes of programs generated by the translation suffer from space inefficiency (16), as a consequence of the transformation making changes to data structures defined in the program. Our transformation involves no such changes.

The NESL (17) and Data Parallel Haskell (42) compilers implement fusion transformation in order to increase granularity. Fusion transforms the program to eliminate redundant synchronization points and intermediate arrays. Although fusion reduces scheduling costs by combining adjacent parallel loops, it is not relevant to controlling granularity within loops. As such, fusion is orthogonal to our oracle based approach.

#### 10.5 Scheduling and Locality

In our approach, we make scheduling decisions based on a property of the intrinsic semantics of pieces of computation: the work of parallel threads. In other related work on scheduling for data locality (3; 15), recent work showed that basing scheduling decisions on the space usage of parallel threads, a different property of intrinsic semantics, can be helpful in improving locality. For example, a class of “space-bounded schedulers” (22; 14; 23) are shown to be well suited for improved locality on deep memory hierarchies. Since space usage of a parallel thread can be bounded by its work, our techniques may be helpful in scheduling for improved data locality.

### **10.6 Cost Semantics**

To give an accurate accounting of thread-creation of overheads in implicitly parallel languages we use a cost semantics, where evaluation steps (derivation rules) are decorated with work and span information or “costs”. This information can then be used directly to bound running time on parallel computers by using standard scheduling theorems that realize Brent’s bound. Many previous approaches also use the same technique to study work-span properties, some of which also make precise the relationship between cost semantics and the standard directed-acyclic-graph models (13; 16; 51). The idea of instrumenting evaluations to generate cost information goes back to the early 90s (48; 46).

### **10.7 Inferring Complexity Bounds**

Our implementation of oracle scheduling requires the programmer to enter complexity bounds for all parallel tasks. In some cases, these bounds can be inferred by various static analyses, for example, using type-based and other static analyses (e.g., (24; 36)), symbolic techniques (e.g., (31; 32)). Our approach can benefit from these approaches by reducing the programmer burden, making it ultimately easier to use the proposed techniques in practice.

### **10.8 Further Validation**

Since the publication of the conference version of this paper (4), we have conducted additional research that further validates the proposed techniques, extends their applicability, and explores extensions to other problem domains outside of the fork-join idiom considered here. In addition to the ML library described in this paper, we implemented our techniques as an optimized, lower-level library for the C++ language (6), which uses a highly optimized work-stealing scheduler to achieve efficiency (5). This library offers techniques for writing parallel programs in the full generality of the C++ language and is now a reasonably mature software system. We used the library in undergraduate teaching at Carnegie Mellon University (1; 2) and elsewhere in intensive courses on parallelism. Our experience in developing this library and the feedback from over 500 students taught thus far suggest that the techniques proposed in this paper can be implemented and used in practice. More recently, we have also started researching the granularity problem in more unstructured parallel-computing problems, including for example graphs, and presented an algorithm for unordered parallel depth-first-search that achieves its efficiency by carefully controlling granularity (7). We are currently working on extending the automatic granularity-control techniques presented here to such unstructured problems.

## **11 Challenges for Programming with Complexity Functions**

The oracle-guided semantics that we presented is relatively general purpose and broadly applicable. Our implementation of this semantics based on the approximation algorithm (Section 7) and the translation (Section 8), however, make the following two assumptions.

- Each complexity function in the given parallel program is accurate and efficient.
- Our approximation algorithm can always approximate constant factors effectively.

Although these two requirements are easily met for large classes of programs, there are challenging cases. In what follows, we describe several such challenging cases that we have encountered.

### ***11.1 Complexity Functions Requiring Auxiliary Data Structures***

There are programs for which providing a constant-time function that computes the complexity of a computation requires the pre-computation of an auxiliary data structure. We ran into such an example when implementing the sparse-matrix by dense vector multiplication benchmark program. Essentially, we needed to pre-compute a prefix sum array over the input data in order to determine the number of non-zero values covered by a range of consecutive rows in the sparse matrix. Fortunately, the cost of this pre-computation turned out to be relatively small in front of the rest of the computations, so the overhead associated with the introduction of the complexity function was limited. For additional details on this example, we refer to the evaluation section of an earlier description of our work (4).

### ***11.2 Complexity Functions for Higher-Order Functions***

Another challenge is that of providing complexity functions for higher-order functions. Consider for example a “map” function that can be used to apply a given function  $f$  to all the leaves of a binary tree. Clearly, the asymptotic complexity of a call of the map function to a function  $f$  and a tree depends not just on the tree but also on the evaluation of the asymptotic complexity function  $f$  on the various leaves of the tree. It is far from obvious how to handle the general case properly. Nevertheless, there are at least two large classes of  $f$  functions for which we are able to implement a useful asymptotic complexity function for predicting the cost of calls to map.

The first class includes constant-time functions. If we map a constant function to all the leaves of a binary tree, then the asymptotic cost is proportional to the number of leaves in the tree. Therefore, to evaluate the cost of the map function on any subtree, it suffices to precompute and store, for each node from the tree, the size of the subtree rooted at this node. Of course, caching additional values in the data structures and maintaining those values upon changes induce additional overheads, but it seems hard to control granularity without any form of additional information.

This approach generalizes to the application of map to functions that do not necessarily take constant time but take a time that depends on the *weight* of the leaf it is applied to, for some definition of *weight*. For example, consider a binary tree that stores lists of integers in its leaves, and assume that we map to this tree a function that increments all the integers stored in the leaves. In this case, we define the weight of a leaf as the length of the list it stores, and define the weight of a subtree as the sum of the weight of its leaves. If we pre-compute the weight of every node in the tree, then we are able to provide a useful asymptotic cost function.

The implementation of the complexity function is not the only challenge associated with higher-order functions: evaluation of the constant factors is another. For example, the constant factors of the map function depend both on the constant factors of the function  $f$  passed to it and on constant factors inherent to the traversal of tree data structures. One

possible solution is to allocate one different estimator for every instantiation of the map function. However, doing so results in imperfect sharing of constant factors, possibly leading to longer convergence phases. In summary, higher-order functions are associated with a number of open challenges.

### *11.3 Mismatch between Average and Worst-Case Complexity*

Our ability to implement the oracle relies on the assumption that we are able to predict asymptotically the amount of work involved in a computation. However, there are cases where the amount of computation to be performed depends on the data itself, and not just on the size of the data.

For example, consider string comparison. On many pairs of input, the comparison terminates quickly because of a mismatch occurring in the first few characters of the two strings. However, the worst-case execution time is  $O(n)$ . If we naively use  $O(n)$  as the complexity function, but the execution is most frequently  $O(1)$ , we will conclude that the constant factors are tiny. At this point, if we receive a pair of two identical strings, we will incorrectly predict the comparison to run fast, and as a result we will fully sequentialize the linear-time comparison of the two strings, even though it could have been performed faster in parallel.

Another example is that of (non randomized) quick sort. The expected complexity is  $O(n \log n)$ , but in the worst case it can run in  $O(n^2)$ . We are tempted to use  $O(n \log n)$  for making predictions. What happens if we do so but then run on an input that triggers a quadratic behavior? We will likely sequentialize sub-computations of much larger size than we ideally aim for. However, it turns out that the inputs that triggers quadratic behavior are also those for which the quick sort algorithm does not exhibit any parallelism. In other words, we could not have sorted the items faster using quick sort by controlling granularity differently.

### *11.4 Initial Estimation of the Constant Factors*

The approach described in this paper has another limitation related to the way that we currently handle the initial convergence phase for the estimation of the constant factors. This limitation may affect code whose structure is more complex than simple recursive divide-and-conquer functions.

To understand the limitation, recall that, when a program begins, we do not have any estimation for its constant factors, and that we initially consider a pessimistic (i.e. large) value for the constants. By doing so, we typically begin by sequentializing only tiny pieces of work; later, when a constant factor is measured to be actually smaller than its initial pessimistic value, we are then able to sequentialize larger pieces of work. This scheme suffices in the case of simple recursive functions, because the constant factors that are measured one level of recursion can be exploited to make predictions at a higher level of recursion.

Consider, however, the case of two nested loops, both implemented as basic divide-and-conquer recursive functions. The constant-factor estimation for the inner loop will converge properly. However, the constant-factor estimation for the outer loop may never get

updated because a measured run is not necessarily triggered. Even when the full execution of the inner loop is systematically sequentialized, consecutive iterations of the outer loop never get sequentialized, because the estimator of the outer loop, which does not share information with the estimator of the inner loop, continues to use the initial pessimistic constant estimation.

Intuitively, our approximation algorithm is facing a bootstrapping problem here: on the one hand, for sequentializing pieces of work, we need some reasonably accurate estimate of the constant factors; on the other hand, we may only obtain these estimations through sequentialized computations. We leave to future work the design of better approximation algorithms and runtime techniques for either propagating more information or for estimating upper bounds to the constant factors, so as to be able to escape this bootstrapping issue.

## 12 Conclusion

In this paper we present an analysis of the impact of thread-creation costs on the performance of implicitly parallel programs and provide a solution for controlling these costs based on an oracle. We formulate our solutions in the context of a nested-parallel functional language and prove that it successfully controls thread-creation costs by using an amortization technique, known also as granularity control. We then present an approximation algorithm for realizing our oracle-guided semantics with the help of asymptotic complexity functions provided by the programmer. We implement our techniques and perform an experimental evaluation, which shows that the approach is effective in controlling the thread-creation costs without detrimentally affecting the performance in the benchmarks considered. While this paper makes important progress towards solving the granularity-control problem, it also shows that there is more work to be done. For example, we have seen theoretically that our approximation algorithm only works for certain kinds of computations. While we have also validated the techniques proposed in this paper in other research projects and in teaching (Section 10.8), our implementations and experiments are still relatively new and would benefit from a more thorough evaluation that considers larger and more varied benchmarks.

## References

- Umut A. Acar and Guy Blelloch. 15210: Algorithms: Parallel and sequential, 2015. <http://www.cs.cmu.edu/~15210/>.
- Umut A. Acar and Guy Blelloch. *Algorithm Design: Parallel and Sequential*. 2015. <http://www.parallel-algorithms-book.com>.
- Umut A. Acar, Guy E. Blelloch, and Robert D. Blumofe. The data locality of work stealing. *Theory of Computing Systems*, 35(3):321–347, 2002.
- Umut A. Acar, Arthur Charguéraud, and Mike Rainey. Oracle scheduling: Controlling granularity in implicitly parallel languages. In *ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*, 2011.
- Umut A. Acar, Arthur Charguéraud, and Mike Rainey. Scheduling parallel programs by work stealing with private dequeues. In *PPoPP '13*, 2013.

- Umut A. Acar, Arthur Charguéraud, and Mike Rainey. *An Introduction to Parallelism in C++*. 2015. Manuscript available at <http://www.cs.cmu.edu/~15210/pas1.html>.
- Umut A. Acar, Arthur Charguéraud, and Mike Rainey. A work-efficient algorithm for parallel unordered depth-first search. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC '15*, pages 67:1–67:12, 2015.
- Gad Aharoni, Dror G. Feitelson, and Amnon Barak. A run-time algorithm for managing the granularity of parallel functional programs. *Journal of Functional Programming*, 2:387–405, 1992.
- Nimar S. Arora, Robert D. Blumofe, and C. Greg Plaxton. Thread scheduling for multiprogrammed multiprocessors. In *Proceedings of the tenth annual ACM symposium on Parallel algorithms and architectures, SPAA '98*, pages 119–129. ACM Press, 1998.
- Nimar S. Arora, Robert D. Blumofe, and C. Greg Plaxton. Thread scheduling for multiprogrammed multiprocessors. *Theory Comput. Syst.*, 34(2):115–144, 2001.
- Josh Barnes and Piet Hut. A hierarchical  $O(N \log N)$  force calculation algorithm. *Nature*, 324:446–449, December 1986.
- Lars Bergstrom, Matthew Fluet, Mike Rainey, John Reppy, and Adam Shaw. Lazy tree splitting. In *ICFP 2010*, pages 93–104. ACM Press, September 2010.
- Guy Blelloch and John Greiner. Parallelism in sequential functional languages. In *Proceedings of the 7th International Conference on Functional Programming Languages and Computer Architecture*, pages 226–237, 1995.
- Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, and Harsha Vardhan Simhadri. Scheduling irregular parallel computations on hierarchical caches. In *Proceedings of the 23rd ACM Symposium on Parallelism in Algorithms and Architectures, SPAA '11*, pages 355–366, 2011.
- Guy E. Blelloch and Phillip B. Gibbons. Effectively sharing a cache among threads. In *SPAA*, 2004.
- Guy E. Blelloch and John Greiner. A provable time and space efficient implementation of NESL. In *Proceedings of the 1st ACM SIGPLAN International Conference on Functional Programming*, pages 213–225. ACM, 1996.
- Guy E. Blelloch, Jonathan C. Hardwick, Jay Sipelstein, Marco Zagha, and Siddhartha Chatterjee. Implementation of a portable nested data-parallel language. *J. Parallel Distrib. Comput.*, 21(1):4–14, 1994.
- Guy E. Blelloch and Gary W. Sabot. Compiling collection-oriented languages onto massively parallel computers. *Journal of Parallel and Distributed Computing*, 8:119–134, February 1990.
- Robert D. Blumofe and Charles E. Leiserson. Scheduling multithreaded computations by work stealing. *J. ACM*, 46:720–748, September 1999.
- Richard P. Brent. The parallel evaluation of general arithmetic expressions. *J. ACM*, 21(2):201–206, 1974.
- Manuel M. T. Chakravarty, Roman Leshchinskiy, Simon Peyton Jones, Gabriele Keller, and Simon Marlow. Data parallel Haskell: a status report. In *Workshop on Declarative Aspects of Multicore Programming, DAMP '07*, pages 10–18, 2007.
- R.A. Chowdhury, F. Silvestri, B. Blakeley, and V. Ramachandran. Oblivious algorithms for multicores and network of processors. In *International Symposium on Parallel Distributed Processing (IPDPS)*, pages 1–12, April 2010.
- Richard Cole and Vijaya Ramachandran. Resource oblivious sorting on multicores. In *ICALP*, 2010.
- Karl Crary and Stephanie Weirich. Resource bound certification. In *Proceedings of the 27th ACM SIGPLAN-SIGACT symposium on Principles of programming languages, POPL '00*, pages 184–198, 2000.
- Marc Feeley. A message passing implementation of lazy task creation. In *Parallel Symbolic Computing*, pages 94–107, 1992.

- Marc Feeley. *An efficient and general implementation of futures on large scale shared-memory multiprocessors*. PhD thesis, Brandeis University, Waltham, MA, USA, 1993. UMI Order No. GAX93-22348.
- Matthew Fluet, Mike Rainey, and John Reppy. A scheduling framework for general-purpose parallel languages. In *ACM SIGPLAN International Conference on Functional Programming (ICFP)*, pages 241–252. ACM, 2008.
- Matthew Fluet, Mike Rainey, John Reppy, and Adam Shaw. Implicitly threaded parallelism in Manticore. *Journal of Functional Programming*, 20(5-6):1–40, 2011.
- Jeremy D. Frens and David S. Wise. Auto-blocking matrix-multiplication or tracking blas3 performance from source code. In *Proceedings of the sixth ACM SIGPLAN symposium on Principles and practice of parallel programming, PPOPP '97*, pages 206–216, New York, NY, USA, 1997. ACM.
- Matteo Frigo, Charles E. Leiserson, and Keith H. Randall. The implementation of the Cilk-5 multithreaded language. In *PLDI*, pages 212–223, 1998.
- Simon F. Goldsmith, Alex S. Aiken, and Daniel S. Wilkerson. Measuring empirical computational complexity. In *Proceedings of the 6th Joint Meeting of the European Software Engineering Conference and the ACM Symposium on the Foundations of Software Engineering*, pages 395–404, 2007.
- Sumit Gulwani, Krishna K. Mehra, and Trishul Chilimbi. Speed: precise and efficient static estimation of program computational complexity. In *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 127–139, 2009.
- Robert H. Halstead. Multilisp: a language for concurrent symbolic computation. *ACM Transactions on Programming Languages and Systems*, 7:501–538, 1985.
- Tasuku Hiraishi, Masahiro Yasugi, Seiji Umatani, and Taiichi Yuasa. Backtracking-based load balancing. In *PPoPP '09*, pages 55–64. ACM, 2009.
- Lorenz Huelsbergen, James R. Larus, and Alexander Aiken. Using the run-time sizes of data structures to guide parallel-thread creation. In *Proceedings of the 1994 ACM conference on LISP and functional programming, LFP '94*, pages 79–90, 1994.
- Steffen Jost, Kevin Hammond, Hans-Wolfgang Loidl, and Martin Hofmann. Static determination of quantitative resource usage for higher-order programs. In *Principles of Programming Languages (POPL)*, pages 223–236, 2010.
- Xavier Leroy, Damien Doligez, Jacques Garrigue, Didier Rémy, and Jérôme Vouillon. *The Objective Caml system*, 2005.
- P. Lopez, M. Hermenegildo, and S. Debray. A methodology for granularity-based control of parallelism in logic programs. *Journal of Symbolic Computation*, 21:715–734, June 1996.
- Eric Mohr, David A. Kranz, and Robert H. Halstead Jr. Lazy task creation: a technique for increasing the granularity of parallel programs. In *Conference record of the 1990 ACM Conference on Lisp and Functional Programming*, pages 185–197, New York, New York, USA, June 1990. ACM Press.
- Girija Jayant Narlikar. *Space-efficient scheduling for parallel, multithreaded computations*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, USA, 1999.
- Joseph Pehoushek and Joseph Weening. Low-cost process creation and dynamic partitioning in Qlisp. In Takayasu Ito and Robert Halstead, editors, *Parallel Lisp: Languages and Systems*, volume 441 of *Lecture Notes in Computer Science*, pages 182–199. Springer Berlin / Heidelberg, 1990.
- Simon L. Peyton Jones. Harnessing the multicores: Nested data parallelism in Haskell. In *APLAS*, page 138, 2008.
- Simon L. Peyton Jones, Roman Leshchinskiy, Gabriele Keller, and Manuel M. T. Chakravarty. Harnessing the multicores: Nested data parallelism in Haskell. In *FSTTCS*, pages 383–414, 2008.

- H. C. Plummer. On the problem of distribution in globular star clusters. *Monthly Notices of the Royal Astronomical Society*, 71:460–470, March 1911.
- Mike Rainey. *Effective Scheduling Techniques for High-Level Parallel Programming Languages*. PhD thesis, University of Chicago, August 2010.
- Mads Rosendahl. Automatic complexity analysis. In *FPCA '89: Functional Programming Languages and Computer Architecture*, pages 144–156. ACM, 1989.
- Daniel Sanchez, Richard M. Yoo, and Christos Kozyrakis. Flexible architectural support for fine-grain scheduling. In *Proceedings of the fifteenth edition of ASPLOS on Architectural support for programming languages and operating systems*, ASPLOS '10, pages 311–322, New York, NY, USA, 2010. ACM.
- David Sands. *Calculi for Time Analysis of Functional Programs*. PhD thesis, University of London, Imperial College, September 1990.
- K. C. Sivaramakrishnan, Lukasz Ziarek, and Suresh Jagannathan. Multimlton: A multicore-aware runtime for standard ml. *Journal of Functional Programming*, FirstView:1–62, 6 2014.
- Daniel Spoonhower. *Scheduling Deterministic Parallel Programs*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, USA, 2009.
- Daniel Spoonhower, Guy E. Blelloch, Robert Harper, and Phillip B. Gibbons. Space profiling for parallel functional programs. In *International Conference on Functional Programming*, 2008.
- Alexandros Tzannes, George C. Caragea, Uzi Vishkin, and Rajeev Barua. Lazy scheduling: A runtime adaptive scheduler for declarative parallelism. *TOPLAS*, 36(3):10:1–10:51, September 2014.
- Leslie G. Valiant. A bridging model for parallel computation. *CACM*, 33:103–111, August 1990.
- Joseph S. Weening. *Parallel Execution of Lisp Programs*. PhD thesis, Stanford University, 1989. Computer Science Technical Report STAN-CS-89-1265.