

Responsive Parallel Computation: Bridging Competitive and Cooperative Threading

Stefan K. Muller

Carnegie Mellon University, USA
smuller@cs.cmu.edu

Umut A. Acar

Carnegie Mellon University, USA;
Inria, France
umut@cs.cmu.edu

Robert Harper

Carnegie Mellon University, USA
rwh@cs.cmu.edu

Abstract

Competitive and cooperative threading are widely used abstractions in computing. In competitive threading, threads are scheduled preemptively with the goal of minimizing response time, usually of interactive applications. In cooperative threading, threads are scheduled non-preemptively with the goal of maximizing throughput or minimizing the completion time, usually in compute-intensive applications, e.g. scientific computing, machine learning and AI.

Although both of these forms of threading rely on the same abstraction of a thread, they have, to date, remained largely separate forms of computing. Motivated by the recent increase in the mainstream use of multicore computers, we propose a threading model that aims to unify competitive and cooperative threading. To this end, we extend the classic graph-based cost model for cooperative threading to allow for competitive threading, and describe how such a cost model may be used in a programming language by presenting a language and a corresponding cost semantics. Finally, we show that the cost model and the semantics are realizable by presenting an operational semantics for the language that specifies the behavior of an implementation, as well as an implementation and a small empirical evaluation.

CCS Concepts • Software and its engineering → Parallel programming languages; Concurrent programming languages; Functional languages

Keywords Parallelism, Cost Models, Cost Semantics, Operational Semantics, Scheduling

1. Introduction

The idea of multiple threads sharing an address space is one of the most widely applicable abstractions in computer

science. Over many years of research and practice, two forms of threading have emerged: competitive threading and cooperative threading. Although they both rely on essentially the same abstraction of threads, these two forms of threading differ and complement each other in their domain of applications, the form of scheduling that they use, and their performance goals, as summarized by the table below.

	Application	Scheduling	Goal
Competitive	Interactive	Preempt.	Responsiveness
Cooperative	Parallel	Non-preemp.	Throughput

Broadly used in interactive systems [33], the work on competitive threads goes back to early systems such as Xerox’s STAR [57] and Cedar [60]. Such systems rely on threads to implement responsive interaction between the different components of the systems (e.g., I/O subsystem, the network) and between the system and the users [33]. Maximizing *responsiveness* is the main performance goal in interactive systems, since this is key to the user experience. To this end, threads are scheduled preemptively, often based on priorities [7, 24, 27, 33].

Cooperative threading is broadly used in fine-grained parallelism, and its use goes back to early parallel programming languages such as Id [5] and Multilisp [30], but it has regained fresh popularity with the increasing mainstream availability of multicore computers. Parallel applications, usually drawn from areas such scientific computing, physical simulations, machine learning and AI, and discrete optimization, are usually compute-intensive and use threads to reduce execution time. To this end, they break up the computation into smaller threads that can be run in parallel and rely on a non-preemptive scheduler to map the threads onto processors. The goal of the scheduler is to minimize the execution time of a parallel application by maximizing *throughput*.

It is technically possible to use competitive threading to implement parallel programs by, for example, creating a small number of system threads and manually scheduling the work of an application over them. This approach, however, can result in complex, low-level, and error-prone code. There has therefore been much work on specialized programming

languages and language extensions for parallel systems, including NESL [10], OpenMP, Cilk [26], Fork/Join Java [41], X10 [18], TBB [36], TPL [42], parallel Haskell [17, 39], parallel ML [25, 37, 51] and Habanero Java [35]. To ensure high performance, these systems rely on non-preemptive schedulers such as work stealing [2, 4, 13], depth-first schedulers [11], and priority schedulers [34].

As shared memory multicore computers have become the common platform for essentially all applications, ranging from compute-intensive to interactive, many applications would benefit from a threading model that bridges competitive and cooperative threading. In such a model, an application can create both competitive and cooperative threads and expect them to be scheduled optimally, that is, to maximize both throughput and responsiveness. For example, an application that interacts with a user as it also performs parallel compute-intensive tasks mixes throughput-oriented parallel computation with responsiveness-oriented interaction.

In this paper, we propose a language and accompanying cost model that combines the competitive and cooperative threading models. We build on a popular graph-based cost model for parallel computing (e.g. [13, 38]), which goes back to the 1960s [28], in which an execution of a parallel program is represented with a Directed Acyclic Graph (DAG or simply dag). We extend this model (Section 2) to allow instructions to be assigned priorities, *foreground* and *background*, which correspond to high-priority and low-priority. We then present a scheduling principle, called *prompt scheduling*, which generalizes the standard *greedy scheduling* [4, 15] to bound both the run-time and responsiveness of a computation. To establish the bound, we make an important assumption that requires the absence of priority inversions in which high-priority computations depend on low-priority ones.

Like all other dag-based cost models, our model allows reasoning about run-time and responsiveness but it leaves an important gap: it only applies to a specific execution of the program rather than the program. To close this gap, we present a small core language (Section 3), called λ^p , which we equip with a cost semantics, following prior language-based cost models [8, 9, 29]. The language supports cooperative threading based on the popular fork-join paradigm and competitive threading based on two constructs for associating priorities with computations. Furthermore, it has a type system based on linear temporal logic that guarantees that well-typed programs avoid the above-mentioned priority inversions. The result is the ability to reason about cost at the level of the program rather than that of the execution.

The dag-based cost models and the cost semantics are both abstract notions of cost that have little value unless they can be realized by an implementation. We show that the cost semantics of λ^p is theoretically *realizable* by giving a transition system (Section 4) that specifies the implementation of a language runtime, and proving that the operational semantics matches the cost semantics. The operational semantics

```
function fib n =
  if n <= 1 then n
  else
    let (a, b) = par (fib (n - 1), fib (n - 2))
    in a + b
```

Figure 1. Code for parallel Fibonacci.

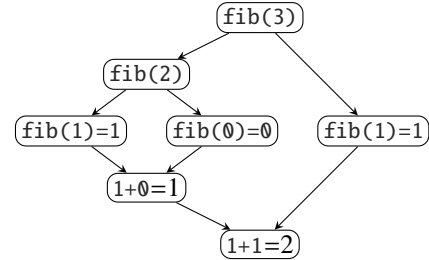


Figure 2. A dag representation of fib(3).

captures important details of an implementation and can be implemented on a modern multicore machine by providing a scheduling algorithm. We briefly describe (Section 5) how such a scheduling algorithm may be implemented.

Finally, we present a prototype implementation of the proposed techniques as an extension to the MLton compiler for Standard ML and perform a small empirical evaluation. Our results show that our theoretical bounds predict the practical run-time and responsiveness of a number of interactive parallel programs.

2. The DAG Model and Prompt Scheduling

The Standard DAG Model. It is common to represent parallel computations using directed acyclic graphs or *dags*. Vertices of the dag represent instructions of the computation, each of which executes in one unit of time, which we call a *step*. Edges represent dependencies between instructions: an edge from u to u' indicates that the instruction represented by u must execute before u' . For a dag g , we write $u \leq_g u'$ to indicate that u is an ancestor of u' in g . When it is clear from the context, we drop g and simply write $u \leq u'$.

For example, consider the function $\text{fib}(n)$, which computes the n^{th} Fibonacci number by performing the two recursive subcalls $\text{fib}(n-1)$ and $\text{fib}(n-2)$ in parallel¹. Figure 1 shows the code for fib . We can represent an execution of $\text{fib}(3)$ as a dag, as shown in Figure 2. For brevity, each vertex represents a call to fib instead of an individual instruction, but can be expanded into a chain of instructions if desired. Vertices with out-degree two “fork” two parallel computations, which may be executed in two (cooperative) threads. Vertices with in-degree two “join” two parallel computations; a join vertex synchronizes its two in-neighbors by waiting for both of them to complete before executing.

¹ While inefficient, this algorithm is commonly used in the literature to illustrate a simple compute-intensive parallel computation.

```

function hello i =
  if i <= 0 then bg ()
  else
    let _ = output("What is your name?")
        x = input ()
        _ = output("Hello, " ^ x)
    in
      hello (i-1)

function fib_hello () = par(fib 3, fg (hello 1))

```

Figure 3. Fibonacci composed with an interactive process.

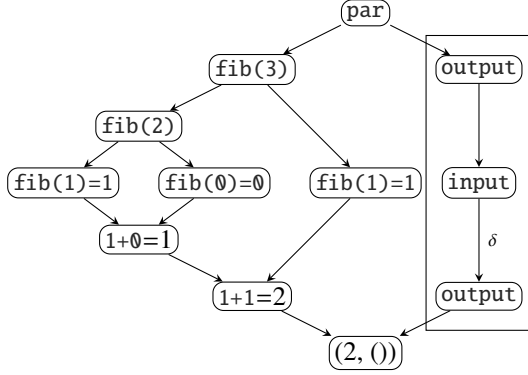


Figure 4. A dag representation of fib_hello.

Our Model. To model responsiveness concerns, we extend the standard dag model to allow certain portions of a dag, called *foreground blocks*, to be specified as *foreground* or high-priority computations. A foreground block is specified by its source and sink vertices. A foreground block with source s and sink t is written $\langle s_t \rangle$ and is the vertex-induced subdag of g consisting of all u such that $s \leq_g u \leq_g t$.

To account for the latency incurred by input operations, we weight edges with the number $\delta \in \mathbb{N}$ of steps by which the operation represented by the source vertex is delayed [45]. More specifically, for a weighted edge (u, u', δ) , if $\delta = 1$, then u incurs no latency and u' may execute on the next step. If $\delta > 1$, then u incurred a latency of δ and u' may execute anytime δ steps after u starts executing.

Mathematically speaking, a dag is a tuple (s, t, V, E, F) consisting of a source vertex s , a sink vertex t , a set V of vertices (where $s, t \in V$ and $s \leq_g t$), a set E of weighted, directed edges, and a set F of foreground blocks. We will derive run-time and responsiveness properties for dags that have no priority inversions, in which a high-priority computation depends on a low-priority one. Without this property, which we call *well-formedness*, we cannot ensure responsiveness. We say that a dag is *well-formed* if each foreground block satisfies the condition that no vertex in the block except for the source has an incoming edge from outside the block. That is, for all $\langle s_t \rangle \in F$ and all $u \neq s \in \langle s_t \rangle$, there does not exist $(u', u, \delta) \in E$ for any $u' \notin \langle s_t \rangle$.

As an example, consider the simple program shown in Figure 3. The function `fib_hello` mixes computation and

interaction by computing `fib(3)` and, in parallel, asking the user a question and responding to the user's answer. The keyword `fg` indicates that the interaction should be given high priority (i.e. is a foreground computation). Figure 4 illustrates the dag for this program. The foreground computation is drawn within a box. The edge weight δ stands for the latency incurred by the input instruction. For all other edges, where the edge weight is 1, we don't explicitly write the weight.

Cost Metrics: Work, Span, Width. In parallel computing with cooperative threads, the *work* of a dag g , which we write $W(g)$, is defined as the number of vertices in the dag and *span* $S(g)$ is defined as the length of the longest path in the dag. When edge weights are used to account for latency, work remains the same as in the traditional model, because time spent blocking on inputs requires delay but no computational work. The span, on the other hand, is now the longest *weighted* path in the dag. The span takes the delays into account since the computation cannot complete until all of the inputs are available [45]. As usual, span corresponds to the time needed to complete the computation with infinitely many processors. Work now corresponds to the total active processing time of the computation. With latencies, it may not be possible to complete the computation in time $W(g)$.

The rest of this section extends the model to account for prioritized computations and uses the extensions to bound both the completion time and the responsiveness of interactive parallel computations. No additional changes are necessary to the notions of work and span beyond including edge weights in the span. However, we distinguish between total and foreground-only work and span. The bounds on the response time will involve the work and span of only the foreground blocks, reflecting the desire that the amount of low-priority computation should not affect responsiveness. For a foreground block f (which, recall, is itself a subdag of the overall dag of the computation), we write $W(f)$ and $S(f)$ for the work and span, respectively of the block. For a graph $g = (s, t, V, E, F)$, the *foreground work* $W^\circ(g)$ and *foreground span* $S^\circ(g)$ are the sum over all foreground blocks:

$$\begin{aligned}
 W^\circ(g) &\triangleq \sum_{f \in F} W(f) \\
 S^\circ(g) &\triangleq \sum_{f \in F} S(f)
 \end{aligned}$$

To bound the response time, we define a new notion, called *foreground width*, which intuitively corresponds to the maximum number of foreground blocks that can be executing at the same time. Formally, we say that two foreground blocks f_1 and f_2 are *serial* if there exists a directed path in the graph from a vertex of f_1 to a vertex of f_2 or vice versa. A set of foreground blocks $F' \subset F$ is *independent* if for all $f_1, f_2 \in F'$, f_1 and f_2 are not serial. The *foreground width* $D(g)$ of a graph g is

$$D(g) \triangleq \max \{ |F'| \mid F' \subset F \wedge F' \text{ is independent} \}$$

Prompt Schedules. A schedule is an assignment of vertices to processors at each step such that if a vertex u is executed

at step i , it is *ready* at step i . A vertex u is ready if all of its ancestors have executed and its latency requirements (if any) have expired. A schedule is *greedy* if as many ready vertices as possible are executed at each step. Greedy schedules suffice to minimize run-time (to within a constant factor of optimal), but not response time of high-priority computations. To reduce response time, we propose a generalization of greedy scheduling which we call *prompt scheduling*. We say that a schedule is *prompt* if it is greedy and it also gives priority to foreground blocks, executing as many foreground vertices as possible at each step (up to the number of ready foreground vertices or the number of processors).

Let T_P denote the time to execute a given parallel computation on P processors using a given schedule. Let $f = \langle \begin{smallmatrix} s \\ t \end{smallmatrix} \rangle$ be a foreground block. Given a schedule, we define the P -processor *response time*, $R_P(f)$ as the number of steps between when s becomes ready and when t is executed (inclusive). We define the total P -processor response time, R_P , as the sum of $R_P(f)$ for all foreground blocks in the dag.

The run-time of a greedy schedule of a dag g is bounded by $\frac{W(g)}{P} + S(g)\frac{P-1}{P}$ [12, 20]. Results such as this are well-studied in the literature, and often attributed to Brent [15], who proved a similar result for “level-by-level” schedules. A similar bound exists for weighted dags such as the ones we use [45], but without the $\frac{P-1}{P}$ factor, since in this case it is possible for all processors to be idle at once, which is not possible in the traditional setting without latencies. We generalize this bound to prompt schedules, taking into account both the run-time and the response time. The intuition behind this proof, and many proofs of Brent-type theorems, is that, by definition, a greedy schedule (all prompt schedules are greedy) will either execute P instructions or execute all ready instructions (an entire “level” of the dag), decreasing the critical path by 1. To show the bound on the response time, we similarly show that, if any foreground blocks are ready, each step decreases the foreground work by P or the foreground span by the number of ready foreground blocks.

Theorem 1. *Consider a parallel computation represented by a well-formed dag g with foreground width D . If this computation is scheduled with a prompt schedule, then $T_P \leq \frac{W(g)}{P} + S(g)$ and $R_P \leq D\frac{W^\circ(g)}{P} + S^\circ(g)$.*

Proof. Since all prompt schedules are greedy, the bound on T_P follows from the run-time bound for weighted dags [45]. We now show the bound on the response time.

We split the total response time into two components R_B (for steps when all processors are **Busy** with foreground work) and R_I (for when some processors are **Idle** or not busy with foreground work), which we will bound separately by visualizing each quantity as a bucket to which tokens are added. The total response time is the total number of tokens in R_B and R_I at the end of the computation. We will also need a bucket W_B to track the “busy” component of the work. At a step i , suppose there are n_i ready foreground vertices which

come from N_i foreground blocks. If $n_i \geq P$, then this is a busy step: place this step’s N_i tokens in bucket R_B and P tokens in W_B . If $n_i < P$, then place N_i tokens in R_I .

Since $N_i \leq D$, for every P tokens placed in W_B , at most D tokens are placed in R_B . So, at any time, $\frac{R_B}{D} \leq \frac{W_B}{P}$.

At the end of the computation, the total number of tokens in the work bucket is at most $W^\circ(g)$ since at busy steps, a prompt schedule will execute only foreground vertices. Thus, at the end of the computation, $R_B \leq D\frac{W^\circ(g)}{P}$.

Now consider a token placed in R_I at step i . This token corresponds to a foreground block f for which at least one vertex is ready at step i . Let g_i be the sub-dag consisting of vertices of f that have not been executed after step i . Extend this in the following way to form a dag g_i^* . All vertices and edges in g_i are also in g_i^* . In addition, for all edges (u, u', δ) where u is in $f \setminus g_i$ and u' is in g_i (that is, u has been executed by the end of step i and u' has not), if u was executed in step $i - j$, add to g_i^* vertices $u_1, \dots, u_{\delta-j-1}$ and edges $(u_1, u_2, 1), \dots, (u_{\delta-j-1}, u, 1)$ (that is, add a chain of length $\delta - j - 1$ before u). Note that because g is well-formed, no vertex of g_i may have edges from outside f in g (except the source of f , but the source must be ready or executed at step i or no vertex of f would be ready), and so the vertices of f that are ready at the start of step $i + 1$ are exactly those vertices that are contained in g_i and have in-degree zero in g_i^* . By the definition of a prompt schedule, it must be the case that all ready vertices of f at step i are executed at step i , and so do not appear in g_{i+1}^* . In addition, for any vertex that is incurring latency at the start of step i and so has a chain before it in g_i^* , the chain is decreased by one vertex in g_{i+1}^* . Together, these facts mean that every vertex in g_i^* with in-degree zero is not present in g_{i+1}^* , and so the longest path in g_{i+1}^* is one shorter than the longest path in g_i^* . Since the longest path in g_0^* , by definition, has length $S(f)$, at most $S(f)$ tokens can be placed in R_I corresponding to f . In total, $R_I \leq S^\circ(g)$. Take the total response time to be $R_B + R_I$. \square

Lower Bounds for Online Scheduling. Given a parallel computation represented by a well-formed dag g , Theorem 1 gives an upper bound on the running time and the responsiveness of a prompt schedule. The run-time bound, like the similar bound for greedy schedules, is within a factor of two of optimal, because $W(g)/P$ and $S(g)$ are both, individually, lower bounds on the computation time. We now show a similar result for the bound on the response time under certain conditions: for the bound, we assume an online scheduling algorithm that has no prior knowledge of the computation dag. Specifically, we show that no matter what decisions the scheduler makes, there exists a dag whose response time is no lower than half of the given bound.

Recall that the response time is the sum over all foreground blocks f of the time taken to execute f . Since $S(f)$ is a lower bound on the time to execute f , $S^\circ(g)$, which is the sum of the spans over all blocks, is a lower bound on response time. Thus, to establish a 2-approximation, it suffices to show that DW°/P

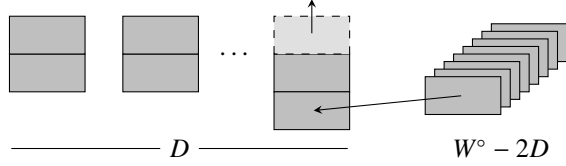


Figure 5. Intuition for the lower bound proof.

is also a lower bound on response time. This is the only part of the argument that relies on the online assumption on the scheduling algorithm. Consider a computation with total work $W^o + 2D$ which consists only of $D \ll W^o$ foreground blocks, each of which is sequential, and the two trees of vertices necessary to fork off and join these foreground blocks. Once all of the foreground blocks have been spawned, think of the work of the computation as W^o “bricks” which are distributed arbitrarily into D stacks, as illustrated in Figure 5. At each step, a prompt scheduler will remove one brick from each of $\min(D, P)$ stacks (foreground blocks). When a stack is empty, that block is complete and no longer counts toward the response time. Since, by assumption, the scheduler only knows which blocks are ready (which stacks have a brick on top) and cannot base its decisions on how large each stack is (this would require knowing how long a block will take to execute), we may play a game against the scheduler. Start by placing two bricks on each stack. Keep the rest of the bricks hidden. At each step, when the scheduler removes a brick from a stack, place another brick at the bottom of that stack until you run out of bricks. In this way, all D blocks will be ready for at least $\frac{W^o - 2D}{\min(D, P)}$ steps (the number of steps it will take to run out of bricks), which will cause the response time to be at least $D \frac{W^o - 2D}{\min(D, P)} \approx D \frac{W^o}{\min(D, P)} \geq D \frac{W^o}{P}$.

3. A Language for Responsive Parallelism

We introduce a core calculus called λ^{ip} , which extends a functional core with constructs for I/O, parallelism and priority. The type system of λ^{ip} separates subcomputations by priority and enforces that high-priority computations do not depend on low-priority ones. The dynamics of λ^{ip} is given by a *cost semantics*, which computes not only the value of an expression, but also an execution dag of the kind described in Section 2, which allows us to reason about cost at the level of the language, and to apply the prompt scheduling theorem to the run-time and responsiveness of programs.

To introduce the features of the language, we consider several simple examples, wherein we use, for convenience, “syntactic sugar,” such as let binding, that can be easily expressed in λ^{ip} . We also make use of base types such as strings and booleans that are not described in the formalism.

3.1 Syntax and Examples

The syntax of λ^{ip} is given in Figure 6. Expressions e include the standard introduction and elimination forms for base types, functions, pairs and sums: natural numbers \mathbf{n} , unit values, λ -abstractions, application, pairs, projection, injec-

<i>Types</i>	$\tau ::=$	$\mathbf{unit} \mid \mathbf{nat} \mid \tau_1 \rightarrow \tau_2 \mid \tau_1 \times \tau_2 \mid \tau_1 + \tau_2 \mid \bigcirc \tau$
<i>Exprs.</i>	$e ::=$	$x \mid \langle \rangle \mid \mathbf{n} \mid \lambda x:\tau.e \mid e e \mid \langle e, e \rangle \mid \mathbf{fst}(e) \mid$ $\mathbf{snd}(e) \mid \mathbf{inl}(e) \mid \mathbf{inr}(e) \mid \mathbf{case}(e)\{x.e; y.e\} \mid$ $\mathbf{fix } x:\tau \mathbf{is } e \mid e \parallel e \mid$ $\mathbf{out}(e) \mid \mathbf{inp}[d](x.e) \mid \mathbf{bg}(e) \mid \mathbf{fg}(e)$

Figure 6. Syntax of λ^{ip} .

tion, and case analysis. We do not include operations (such as $+$ and $<$) on natural numbers for simplicity, but these could be added in a straightforward way. The fixed point operator $\mathbf{fix } x:\tau \mathbf{is } e$ allows expressing recursion. *Parallel tuples* $e_1 \parallel e_2$ (written $\mathbf{par}(e_1, e_2)$ in examples) allow for fork-join parallelism: the expressions e_1 and e_2 denote parallel expressions that may be evaluated in parallel.

Input and Output. The construct $\mathbf{inp}[d](x.e)$ binds user input to the variable x in evaluating e and $\mathbf{out}(e)$ outputs the value of e to the user. The annotation d relates to the cost semantics; we ignore it for now. Our techniques do not make assumptions about how exactly input/output is performed (e.g., via a console, through GUI operations, over a network). We therefore leave these details unspecified for simplicity. Because natural numbers are the only interesting base type of λ^{ip} , only natural numbers can be input/output; generalization to other base types such as strings, as used in our examples, is straightforward. Figure 3 shows an example interactive function `hello` that asks the user questions, repeating for a number of times specified by the argument `i`.

Prioritized Computation. Now consider a parallel interactive program combining `hello` and `fib` from Figure 3:

```
function fib_hello () = par(fib 43, hello 15)
```

This code cannot guarantee responsiveness because it does not distinguish between the competitive thread, executing `hello 15`, and the many low-priority computation threads created by `fib 43`. A scheduler might get lucky, but in general, responses to the user could get arbitrarily delayed as the computation threads starve the interaction.²

As the `fib_hello` example illustrates, we would like to enable the programmer to distinguish between high-priority and low-priority computations. To this end, λ^{ip} provides two language constructs, $\mathbf{fg}(e)$ and $\mathbf{bg}(e)$, that represent, respectively, a *foreground computation* that runs with high priority, and a *background computation* that runs with low priority from within a foreground block. Using these constructs, we can write `fib_hello` so that it runs `hello` in the foreground as shown in Figure 3.

In the example `fib_hello`, foreground and background computations do not interact in interesting ways. For an example where they do, consider a “Fibonacci server”, `fib_server`, shown in Figure 7 on the left. The function asks the user

² Although we do not discuss the details in the paper, we confirmed that indeed such an implementation has poor responsiveness.

```

function fib_server () = function fib_server () =
  let n = input () in    let n = input () in
    if n < 0 then ()    if n < 0 then bg ()
    else                else
      output (fib n);    bg (output (fib n));
      fib_server ()      fib_server ()

function main () =      function main () =
  fib_server ()          fg (fib_server ())

```

Figure 7. Server (l) without priorities and (r) with priorities.

for an input n (a natural number) and computes the n^{th} Fibonacci number using `fib`. Because a Fibonacci computation performs a large amount of work, the input loop could become sluggish. In λ^{ip} , the programmer can solve this problem by running `fib_server` in the foreground and `fib` in the background, as shown on the right in Figure 7. The expression `bg (output (fib n))` spawns a new background thread to perform the Fibonacci computation asynchronously and output the result. The foreground computation can spawn many background computations, each of which computes the requested Fibonacci number in parallel with other background computations as well as the foreground interactive server loop.

3.2 Type System

As presented so far, the language allows “priority inversions” in which foreground code blocks on background computation. As an example, consider the following variant of the Fibonacci server:

```

1 function fib_server_bad () =
2   let n = input () in
3     if n < 0 then bg ()
4     else let fibn = bg (fib n) in
5           output (fg (fibn));
6           fib_server_bad ()
7
8 function main () = fg (fib_server_bad ())

```

The function `fib_server_bad` receives the input n from the user and then creates a background computation `fibn` to compute the n^{th} Fibonacci number (Line 4). It then immediately demands the result for output (Line 5). This program might not be responsive because a foreground computation (function `fib_server_bad`) is waiting on a potentially long-running background computation.

To prevent such responsiveness problems, the type system of λ^{ip} enforces a clean separation between foreground and background code, using techniques inspired by prior type systems for staged computation (Section 6). This separation is sufficient to show (in Section 3.3) that the dags corresponding to well-typed λ^{ip} programs are well-formed in the sense of Section 2 and thus, by Theorem 1, admit prompt schedules that bound responsiveness and completion time. In this section, we describe the salient aspects of the type system.

$$\begin{array}{c}
\frac{}{\Gamma, x : \tau @ w \vdash x : \tau @ w} \quad \frac{}{\Gamma, x : \text{nat} @ w \vdash x : \tau @ w'} \\
\frac{\Gamma, x : \tau @ w \vdash e : \tau' @ w}{\Gamma \vdash \lambda x : \tau . e : \tau \rightarrow \tau' @ w} \quad \frac{\Gamma \vdash e_1 : \tau \rightarrow \tau' @ w \quad \Gamma \vdash e_2 : \tau @ w}{\Gamma \vdash e_1 e_2 : \tau' @ w} \\
\frac{\Gamma \vdash e : \text{nat} @ w}{\Gamma \vdash \text{out}(e) : \text{unit} @ w} \quad \frac{\Gamma, x : \text{nat} @ w \vdash e : \tau @ w}{\Gamma \vdash \text{inp}[d](x.e) : \tau @ w} \\
\frac{\Gamma \vdash e : \tau @ \mathbb{B}}{\Gamma \vdash \text{bg}(e) : \circ \tau @ \mathbb{F}} \quad \frac{\Gamma \vdash e : \circ \tau @ \mathbb{F}}{\Gamma \vdash \text{fg}(e) : \tau @ \mathbb{B}}
\end{array}$$

Figure 8. Selected typing rules for λ^{ip}

The types τ include unit and natural numbers as base types, as well as functions, binary tuples and binary sums and the circle type $\circ \tau$, which represents background computations. The typing judgment has the form $\Gamma \vdash e : \tau @ w$ indicating that e has type τ at “world” w . The world is either \mathbb{F} or \mathbb{B} , indicating that the expression is suitable for the foreground or background, respectively. Contexts Γ have entries of the form $x : \tau @ w$, indicating that variable x is in the context with type τ at world w .

Most of the rules allow expressions to type at any world, but require all subexpressions to be at the same world as the whole expression. Figure 8 shows the typing rules for function types as an example. Most of the rules are similar to these and are omitted for space reasons. Transitions between worlds are effected by the `bg(e)` and `fg(e)` operations. If e has type τ in the background (world \mathbb{B}), then the expression `bg(e)` has the type $\circ \tau$ in the foreground (world \mathbb{F}). This allows encapsulated background computations to be created and passed around in the foreground. If e has type $\circ \tau$ in world \mathbb{F} , then the expression `fg(e)` has type τ in \mathbb{B} . This means that the result of an encapsulated computation can *only* be demanded in the background, which precludes priority inversions. For example, this restriction will rule out the function `fib_server_bad` above, since this function is called in the foreground and the expression `fg fibn` cannot be assigned a type in the foreground.

There are two rules for typing variables. If $x : \tau @ w$ is in the context, the variable x has type τ at world w . We also allow variables of type `nat` to type at either world, allowing foreground code to make use of variables (of type `nat`) bound in the background and vice versa. The restriction to type `nat` ensures that code can’t “escape” to the wrong world encapsulated in a function or thread. This is related to the mobility restriction of Murphy et al. [47], and could easily be expanded to allow any “mobile” type, including sums and products (but not functions or encapsulations of type $\circ \tau$).

3.3 Cost Semantics

We now define a *cost semantics* for λ^{ip} , which both computes the value of an expression and determines an execution dag

of the kind described in Section 2 for a λ^{ip} program. The parallel structure of the program, as well as the cost metrics such as work and span, can be read off from the resulting dag, and are used to reason about the run-time and responsiveness of parallel programs.

The cost semantics is given in Figure 9. The judgment $e \Downarrow^\Delta v; g$ states that the expression e evaluates to v and has cost graph g . The judgment is parametrized by $\Delta : \text{InputIDs} \rightarrow 2^{\mathbb{N}}$, a mapping which assigns a set of possible delays to each input identifier d (recall from the syntax that input operations are tagged with such identifiers). Values v consist of the unit value, numerals, lambda abstractions, pairs and injections of values, and a new form of thread handle which abstractly represents a thread as the value to which it will evaluate and a handle to the sink of its expression’s cost graph:

$$v ::= \langle \rangle \mid \mathbf{n} \mid \lambda x:\tau.e \mid \langle v, v \rangle \mid \mathbf{inl}(v) \mid \mathbf{inr}(v) \mid \mathbf{thread}[u](v)$$

Many of the rules for the sequential components of the language and parallel tuples are based on the cost semantics of Spoonhower et al. [59]. The rules for generating and joining with background threads ($\mathbf{bg}(e)$ and $\mathbf{fg}(e)$, respectively), are based on Spoonhower’s treatment of futures [58], which share the property that an asynchronous expression is spawned in one part of a computation and demanded in another.

The generation of cost graphs is defined as part of the derivation of the evaluation judgment. A graph may consist of a single vertex, written $[u]$, or a single edge, written $[(u_1, u_2, \delta)]$, or may be formed by combining smaller graphs, which are generally produced from evaluating subexpressions. In most cases, subexpressions are evaluated sequentially, represented in the cost graph by combining the cost graphs of the subexpressions using serial composition $g_1 \oplus g_2$ which joins the sink of g_1 to the source of g_2 by an edge of weight 1 (a more general form, \oplus_δ , uses an edge of weight δ , as shown in Figure 11). The empty graph \emptyset is an identity for \oplus . In the rule for $e_1 \parallel e_2$, the cost graphs for e_1 and e_2 are combined using parallel composition $g_1 \otimes g_2$, which joins the graphs in parallel with new vertices s and t as the source and sink (Figure 11). If one of the graphs is empty, the other is simply composed with s and t .

The rule for $\mathbf{bg}(e)$ uses the left parallel composition operator [58]. The graph $g \circledast u$ “hangs g off of” vertex u (Figure 11). For the purposes of sequentially composing this graph with other graphs, u is both the source and the sink, reflecting the fact that the new thread is executed concurrently with the continuation of the current thread.

The rule for $\mathbf{fg}(e)$ evaluates e to a background thread and also gets a handle to the sink of the cost graph for the thread’s expression. The rule adds an edge between the sink and the vertex representing the \mathbf{fg} instruction. In the rule for $\mathbf{fg}(e)$, the cost graph for e is marked as foreground with the operation $g \odot$. This operation produces a foreground block $\langle \overset{s}{\cdot} \rangle$ where s and t are the source and sink of g . Finally,

the input rule adds an edge of weight δ , where δ is chosen nondeterministically from $\Delta(d)$.

Figure 10 formally defines the left parallel composition and foreground block formation rules on graphs. Sequential and parallel composition are omitted for space reasons.

Recall that, in order to apply the results of Section 2 to the dags generated by the cost semantics, we need to show that such dags are well-formed. The well-formedness assumption requires that there are no edges to internal nodes of foreground blocks. Such an edge would correspond to a priority inversion in the language, and is ruled out by the type system, as we will now show. We first show that an expression that types in the foreground will correspond to a dag with no nested foreground blocks or external dependencies.

Lemma 1. *If $\cdot \vdash e : \tau @ \mathbb{F}$ and $e \Downarrow^\Delta v; (s, t, V, E, F)$, then $F = \emptyset$ and for all $(u', u, \delta) \in E$, we have $u \in V$.*

Proof. By induction on the derivation of $\cdot \vdash e : \tau @ \mathbb{F}$. \square

This result can then be easily extended to show that well-typed programs produce well-formed dags.

Theorem 2. *If $\cdot \vdash e : \tau @ w$ and $e \Downarrow^\Delta v; g$, g is well-formed.*

Proof. Let $g = (s, t, V, E, F)$. Proceed by induction on the derivation of $e \Downarrow^\Delta v; g$. The interesting case is the rule for $\mathbf{fg}(e')$, which adds a foreground block. By inversion, $e' \Downarrow^\Delta v'; (s', t', V', E', F')$ and by inversion on the typing rules, $\cdot \vdash e' : \odot \tau @ \mathbb{F}$. By Lemma 1, $F' = \emptyset$ and for all $(u', u, \delta) \in E'$, we have $u' \in V'$. By the cost semantics, we have $F = \{\langle \overset{s'}{\cdot} \rangle\}$ and $E = E' \cup \{(t, u_2, 1), (u_1, u_2, 1)\}$. Since no edge is added with a target in $\langle \overset{s'}{\cdot} \rangle$, there is no $u \in \langle \overset{s'}{\cdot} \rangle$ such that $(u', u, \delta) \in E$ for $u' \notin \langle \overset{s'}{\cdot} \rangle$. No other rule adds an edge to a vertex of a subdag except to its source, so well-formedness is preserved. \square

4. Semantic Realization

We have thus far established bounds on the responsiveness and run-time of prompt schedules of well-formed execution dags (Theorem 1), and defined a language for prioritized interactive parallelism whose cost semantics generates only such dags. The cost model provides a theory of the responsiveness and efficiency of λ^{ip} programs with which we can derive results about programs, but these results remain abstract until we validate them with respect to a lower-level model. In this section, we give a transition semantics that specifies an implementation of λ^{ip} , and show that the cost attributed to a program by the cost model corresponds to a more concrete notion of cost in terms of steps of the transition system. This section will give the broad ideas of the operational semantics and the correspondence proofs, but many details are omitted for space reasons. A full treatment is available in the companion technical report [46].

$$\begin{array}{c}
\frac{}{v \Downarrow^\Delta v; \emptyset} \quad \frac{e_1 \Downarrow^\Delta \lambda x:\tau.e; g_1 \quad e_2 \Downarrow^\Delta v; g_2 \quad [v/x]e \Downarrow^\Delta v'; g_3 \quad u \text{ fresh}}{e_1 e_2 \Downarrow^\Delta v'; g_1 \oplus g_2 \oplus [u] \oplus g_3} \quad \frac{e \Downarrow^\Delta \langle v_1, v_2 \rangle; g \quad u \text{ fresh}}{\text{fst}(e) \Downarrow^\Delta v_1; g \oplus [u]} \\
\frac{e \Downarrow^\Delta \langle v_1, v_2 \rangle; g \quad u \text{ fresh}}{\text{snd}(e) \Downarrow^\Delta v_2; g \oplus [u]} \quad \frac{e_1 \Downarrow^\Delta v_1; g_1 \quad e_2 \Downarrow^\Delta v_2; g_2}{\langle e_1, e_2 \rangle \Downarrow^\Delta \langle v_1, v_2 \rangle; g_1 \oplus g_2} \quad \frac{e_1 \Downarrow^\Delta v_1; g_1 \quad e_2 \Downarrow^\Delta v_2; g_2}{e_1 \parallel e_2 \Downarrow^\Delta \langle v_1, v_2 \rangle; g_1 \otimes g_2} \\
\frac{e \Downarrow^\Delta \text{inl}(v); g_1 \quad [v/x]e_1 \Downarrow^\Delta v'; g_2 \quad u \text{ fresh}}{\text{case}(e)\{x.e_1; y.e_2\} \Downarrow^\Delta v'; g_1 \oplus [u] \oplus g_2} \quad \frac{e \Downarrow^\Delta \text{inr}(v); g_1 \quad [v/y]e_2 \Downarrow^\Delta v'; g_2 \quad u \text{ fresh}}{\text{case}(e)\{x.e_1; y.e_2\} \Downarrow^\Delta v'; g_1 \oplus [u] \oplus g_2} \\
\frac{e \Downarrow^\Delta v; g \quad g = (s, t, V, E, F) \quad u \text{ fresh}}{\text{bg}(e) \Downarrow^\Delta \text{thread}[t](v); g \mathcal{O}^u} \quad \frac{e \Downarrow^\Delta \text{thread}[u_1](v); g \quad u_2 \text{ fresh}}{\text{fg}(e) \Downarrow^\Delta v; (g \mathcal{O}) \oplus [u_2] \cup \{(u_1, u_2, 1)\}} \quad \frac{e \Downarrow^\Delta v; g \quad u \text{ fresh}}{\text{out}(e) \Downarrow^\Delta \langle \rangle; g \oplus [u]} \\
\frac{[n/x]e \Downarrow^\Delta v; g \quad u_1 \text{ fresh} \quad u_2 \text{ fresh} \quad \delta \in \Delta(d)}{\text{inp}[d](x.e) \Downarrow^\Delta v; [u_1] \oplus_\delta [u_2] \oplus g} \quad \frac{[\text{fix } x:\tau \text{ is } e/x]e \Downarrow^\Delta v; g \quad u \text{ fresh}}{\text{fix } x:\tau \text{ is } e \Downarrow^\Delta v; [u] \oplus g}
\end{array}$$

Figure 9. Cost Semantics.

$$\begin{array}{lcl}
[u] & = & (u, u, \{u\}, \emptyset, \emptyset) \\
[(u_1, u_2, \delta)] & = & (u_1, u_2, \{u_1, u_2\}, \{(u_1, u_2, \delta)\}, \emptyset) \\
(s, t, V, E, F) \mathcal{O}^u & = & (u, u, V \cup \{u\}, E \cup \{(u, s, 1)\}, F) \\
(s, t, V, E, F) \mathcal{O} & = & (s, t, V, E, F \cup \{\langle \cdot \rangle\})
\end{array}$$

Figure 10. Selected graph operations.

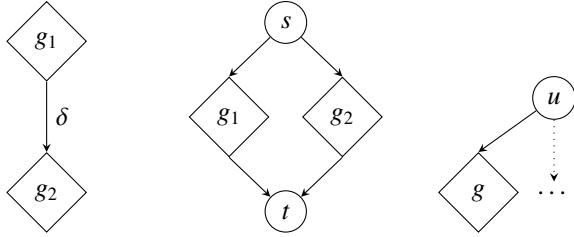


Figure 11. From left: $g_1 \oplus_\delta g_2$, $g_1 \otimes g_2$, and $g \mathcal{O}^u$.

4.1 Operational Semantics

Since the operational semantics is a transition system, it must be able to represent intermediate states of computation. In particular, such an intermediate state may have many active threads of execution. We will name these threads with thread symbols, for which we use the metavariables a, b, c and variants. We use these symbols for threads generated by parallel pairs $e_1 \parallel e_2$ as well as background threads. We also introduce three new expression forms which are not needed for source programs (i.e. programs that have not begun evaluation):

$$e ::= \dots \mid \text{tid}[a] \mid \text{join}[a, b] \mid \text{in}(x.e)$$

The first, $\text{tid}[a]$, is a runtime representation of a background thread identified by thread symbol a . The second, $\text{join}[a, b]$, is a parallel tuple whose components are being evaluated by threads a and b . Finally, $\text{in}(x.e)$ will be used to represent an input expression whose latency has expired but which has not yet produced an input value.

We also introduce *thread pools*. A thread pool μ is a mapping from thread identifiers a to pairs (δ, e) of a delay and an expression, indicating that thread a may run command e after δ steps. We write a thread pool as

$$a_1 \hookrightarrow (\delta_1, e_1) \uplus \dots \uplus a_n \hookrightarrow (\delta_n, e_n)$$

and the concatenation of two disjoint thread pools as $\mu_1 \uplus \mu_2$.

It is straightforward to extend the type system to account for threads. We sketch the extension here. The expression typing judgment becomes $\Gamma \vdash_\Sigma e : \tau @ w$, which includes a thread signature Σ . Thread signatures have entries of the form $a \sim \tau @ w$, indicating that thread a is running an expression of type τ at world w . Most rules are unchanged from Figure 8 and simply pass the thread signature through.

A new typing judgment, $\Gamma \vdash_\Sigma \mu : \Sigma$, indicates that the thread pool μ has the signature Σ . The rules require that $a \sim \tau @ w \in \Sigma$ if and only if $a \hookrightarrow (\delta, e) \in \mu$ and e has type τ at world w .

The operational semantics of λ^p consists of two components: local and global [32]. The local semantics concerns individual threads, and indicates how expressions transition. Selected rules are presented in Figure 12. The rules in this figure correspond to two judgments. The judgment $e \text{ val}$ indicates that e is an irreducible value. Values are the unit value, numerals, functions, pairs and injections of values, and thread handles $\text{tid}[b]$. The local transition judgment is

$$e \mid \mu \mapsto_a^\Delta (\delta', e') \mid \mu \uplus \mu'$$

which states that thread a running e transitions to e' , possibly spawning new threads, which are collected in μ' . The original thread pool μ is unchanged; threads are never altered or removed by local transitions. The thread identifier a is not important for the local transition, but will be used in some of the global definitions and results. The new expression e' will be able to run after a delay of δ' steps (if $\delta' = 0$, it can run immediately). As with the cost semantics, the judgment is parametrized by a delay assignment Δ .

$$\begin{array}{c}
\frac{}{\langle \rangle \text{ val}} \quad \frac{}{\text{n val}} \quad \frac{}{\lambda x:\tau.e \text{ val}} \quad \frac{e_1 \text{ val} \quad e_2 \text{ val}}{\langle e_1, e_2 \rangle \text{ val}} \quad \frac{e \text{ val}}{\text{inl}(e) \text{ val}} \quad \frac{e \text{ val}}{\text{inr}(e) \text{ val}} \quad \frac{}{\text{tid}[a] \text{ val}} \\
\frac{e_1 \mid \mu \mapsto_a^\Delta (\delta, e'_1) \mid \mu \uplus \mu'}{e_1 \ e_2 \mid \mu \mapsto_a^\Delta (\delta, e'_1 \ e_2) \mid \mu \uplus \mu'} \quad \frac{e_2 \mid \mu \mapsto_a^\Delta (\delta, e'_2) \mid \mu \uplus \mu'}{(\lambda x:\tau.e_1) \ e_2 \mid \mu \mapsto_a^\Delta (\delta, (\lambda x:\tau.e_1) \ e'_2) \mid \mu \uplus \mu'} \quad \frac{e_2 \text{ val}}{(\lambda x:\tau.e_1) \ e_2 \mid \mu \mapsto_a^\Delta (0, [e_2/x]e_1) \mid \mu} \\
\frac{b \text{ fresh} \quad c \text{ fresh}}{e_1 \parallel e_2 \mid \mu \mapsto_a^\Delta (0, \text{join}[b, c]) \mid \mu \uplus b \hookrightarrow (0, e_1) \uplus c \hookrightarrow (0, e_2)} \quad \frac{\mu = b \hookrightarrow (\delta_b, e_b) \uplus c \hookrightarrow (\delta_c, e_c) \uplus \mu' \quad e_b \text{ val} \quad e_c \text{ val}}{\text{join}[b, c] \mid \mu \mapsto_a^\Delta (0, \langle e_b, e_c \rangle) \mid \mu} \\
\frac{b \text{ fresh}}{\text{bg}(e) \mid \mu \mapsto_a^\Delta (0, \text{tid}[b]) \mid \mu \uplus b \hookrightarrow (0, e)} \quad \frac{e \mid \mu \mapsto_a^\Delta (\delta, e') \mid \mu \uplus \mu'}{\text{fg}(e) \mid \mu \mapsto_a^\Delta (\delta, \text{fg}(e')) \mid \mu \uplus \mu'} \quad \frac{\mu = b \hookrightarrow (\delta, e) \uplus \mu' \quad e \text{ val}}{\text{fg}(\text{tid}[b]) \mid \mu \mapsto_a^\Delta (0, e) \mid \mu} \\
\frac{e \mid \mu \mapsto_a^\Delta (\delta, e') \mid \mu \uplus \mu'}{\text{out}(e) \mid \mu \mapsto_a^\Delta (\delta, \text{out}(e')) \mid \mu \uplus \mu'} \quad \frac{e \text{ val}}{\text{out}(e) \mid \mu \mapsto_a^\Delta (0, \langle \rangle) \mid \mu} \quad \frac{\delta \in \Delta(d)}{\text{inp}[d](x.e) \mid \mu \mapsto_a^\Delta (\delta - 1, \text{in}(x.e)) \mid \mu} \quad \frac{}{\text{in}(x.e) \mid \mu \mapsto_a^\Delta (0, [n/x]e) \mid \mu}
\end{array}$$

Figure 12. Selected local dynamic rules.

Most of the transition rules are straightforward and are omitted. The complete rules for function application are given as an example: in $e_1 \ e_2$, the subexpression e_1 is stepped until it is a lambda abstraction, then e_2 is stepped until it is a value, which is then substituted for the variable in the body of the abstraction using standard capture-avoiding substitution. A parallel tuple $e_1 \parallel e_2$ spawns two new threads b and c to execute e_1 and e_2 , respectively. The local thread a steps to $\text{join}[b, c]$, indicating that this thread is now waiting for b and c to complete. When both threads have stepped to irreducible values, $\text{join}[b, c]$ steps to a pair of the two values. In the same vein, $\text{bg}(e)$ spawns a new thread b to evaluate e and returns the thread handle $\text{tid}[b]$. Note that, while threads spawned by parallel tuples and threads spawned by $\text{bg}(e)$ are treated identically by the semantics (i.e. they are stepped with the same transitions and not distinguished in the thread pool), the threads b and c spawned by a parallel tuple are never referred to by thread handles (e.g. $\text{tid}[b]$) because these threads are not first class.

The expression $\text{fg}(e)$ steps e until it reaches $\text{fg}(\text{tid}[b])$, which then blocks until thread b has evaluated its expression down to an irreducible value e' , at which point $\text{fg}(\text{tid}[b])$ steps to e' . The input rule is the only one which results in a delay, which is chosen nondeterministically from $\Delta(d)$. After the delay, the new expression $\text{in}(x.e)$ nondeterministically chooses a natural number n to substitute for x in e , representing the uncertainty in the input from the user or environment.

The global rules in Figure 13 define the transitions of entire thread pools, i.e. the entire state of the computation. The judgment $\mu \text{ final}$ states that μ has completed evaluating and its rules simply require that all threads in μ be irreducible. The global step relation is

$$r; \mu \mapsto_{\text{glo}} r'; \mu'$$

and has only one rule, which allows some number N of threads whose delay is 0 to step using the local dynamics.

There is also a counter for the total response time r , which at each step is incremented by the number of ready foreground blocks³ (the formal definition of $RFB(\mu)$, which counts the ready foreground blocks in a thread pool, is straightforward and omitted for simplicity). The new thread pool consists of the updated threads 1 through N , and the unaltered threads $N+1$ through n with their delays (if nonzero) decremented.

Note that the global step relation does not specify a scheduling strategy, nor does it enforce any constraints on schedules other than that only ready threads may step. In our results, we will quantify over valid, prompt schedules: those that step as many threads as possible, prioritizing threads that are executing foreground blocks, bounded by the number of available processors.

We can prove modified progress and preservation results for the λ^p semantics. In addition to type safety, it will become important to show that another property of programs, which we call “well-joinedness”, is preserved throughout execution. Intuitively, well-joinedness (denoted by the judgment $e \text{ wj}$) requires that join expressions appear only in the part of an expression which is currently being evaluated⁴. In particular, they may not appear encapsulated in functions, or in expressions which have not yet been evaluated. Details of the definition of well-joinedness and the safety results are omitted for space reasons.

4.2 Extended Cost Model

In order to show a correspondence between the cost semantics and the operational semantics, we must extend the cost semantics to generate cost graphs not just for expressions

³ Commuting the summations, counting the number of blocks at each step is equivalent to counting the number of steps taken to execute each block, which is the response time.

⁴ For those familiar with evaluation contexts or stack machine semantics, join can only appear in the “hole” of an evaluation context or at the top of a stack.

$$\begin{array}{c}
\frac{}{\emptyset \text{ final}} \quad \frac{e \text{ val} \quad \mu \text{ final}}{a \hookrightarrow (\delta, e) \uplus \mu \text{ final}} \\
\frac{\mu = a_1 \hookrightarrow (\delta_1, e_1) \uplus \dots \uplus a_n \hookrightarrow (\delta_n, e_n) \quad \forall N < i \leq n. \delta'_i = \max(0, \delta_i - 1) \quad N \leq n \quad \forall 1 \leq i \leq N. \delta_i = 0 \quad \forall 1 \leq i \leq N. e_i \mid \mu \mapsto_{a_i}^{\Delta} (\delta'_i, e'_i) \mid \mu \uplus \mu'_i}{r; \mu \mapsto_{\text{glo}} r + |\text{RFB}(\mu)|; a_1 \hookrightarrow (\delta'_1, e'_1) \uplus \dots \uplus a_N \hookrightarrow (\delta'_N, e'_N) \uplus a_{N+1} \hookrightarrow (\delta'_{N+1}, e_{N+1}) \uplus \dots \uplus a_n \hookrightarrow (\delta'_n, e_n) \uplus \mu'_1 \uplus \dots \uplus \mu'_N}
\end{array}$$

Figure 13. Global Dynamics.

Expression cost semantics $e; \mu \Downarrow^{\Delta} v; g$

$$\begin{array}{c}
\frac{\mu = \mu' \uplus b \hookrightarrow (\delta_b, e_b) \uplus c \hookrightarrow (\delta_c, e_c) \quad e_b; \mu \Downarrow^{\Delta} v_1; g_1 \quad e_c; \mu \Downarrow^{\Delta} v_2; g_2 \quad u \text{ fresh}}{\text{join}[b, c]; \mu \Downarrow^{\Delta} \langle v_1, v_2 \rangle; (u, u, \{u\}, \{(b, u, 1), (c, u, 1)\}, \emptyset)} \\
\frac{e; \mu \Downarrow^{\Delta} \text{thread}[u_1](v); g \quad u_2 \text{ fresh}}{\text{fg}(e); \mu \Downarrow^{\Delta} v; (g \odot) \oplus [u_2] \cup \{(u_1, u_2, 1)\}} \\
\frac{\mu = \mu' \uplus b \hookrightarrow (\delta, e_b) \quad e; \mu \Downarrow^{\Delta} \text{tid}[b]; g \quad e_b; \mu \Downarrow^{\Delta} v; g_b \quad u \text{ fresh}}{\text{fg}(e); \mu \Downarrow^{\Delta} v; (g \odot) \oplus [u] \cup \{(b, u, 1)\}}
\end{array}$$

Thread pool cost semantics $\mu_l; \mu_g \Downarrow^{\Delta} v; g$

$$\frac{\mu_l; \mu_g \Downarrow_{\mathcal{C}}^{\Delta} \{G\} \quad e; \mu_g \Downarrow^{\Delta} g \quad g \neq \emptyset}{\emptyset; \mu_g \Downarrow_{\mathcal{C}}^{\Delta} \{ \} \quad a \hookrightarrow (\delta, e) \uplus \mu_l; \mu_g \Downarrow_{\mathcal{C}}^{\Delta} \{a \hookrightarrow g \mid_{\delta} \uplus G\}}$$

Figure 14. Selected extended cost semantics rules.

$$\begin{array}{lcl}
g \Vdash_0 & = & g \\
(s, t, V, E, F) \mid_{\delta} & = & [\alpha] \oplus_{\delta} g \quad \delta > 0, \alpha \text{ fresh} \\
(s, t, V, E, F) \odot & = & (s, t, V, E, F \cup \{\langle s \rangle\}) \quad \nexists a, \delta. (a, s, \delta) \in E \\
(s, t, V, E, F) \ominus & = & (s, t, V, E, F \cup \{\langle t \rangle\}) \quad \exists a, \delta. (a, s, \delta) \in E
\end{array}$$

Figure 15. Extended graph operations.

but also for thread pools which can represent programs that have already begun to execute. As such, dags may no longer have a single source vertex, though they will continue to have a single sink vertex (the final instruction of the initial thread). They will have a source vertex for each ready thread. This modification is relatively straightforward: for each thread, we will generate a standard dag like those of Section 3.3, which we now call a *thread graph* or *thread dag*, with a single source and single sink. These are then composed to form a *configuration graph* or *configuration dag* by adding edges that correspond to the inter-thread dependencies created by `join` and `fg`.

The judgment $e; \mu \Downarrow^{\Delta} v; g$ indicates that the expression e evaluates to v and has cost graph g in the presence of μ . The expression being evaluated may refer to threads in μ . These threads are included so that the value can be generated, but

their cost is not included in g . Most rules for this judgment do not inspect μ and so are similar to the corresponding rules in Figure 9. The rules that are new or substantially different are shown in Figure 14. In addition to the rule for `join`, there is now an additional rule for `fg` handling the case in which e evaluates to a thread handle `tid`[b]. In this case, the expression corresponding to b in the thread pool is evaluated. The definition of $g \odot$ is extended for the case in which the source of g is a join point, i.e. has incoming edges from outside g . This is handled using a new form of foreground block $\langle \cdot \rangle_t$, which has only a sink t and no source. All vertices that are ancestors of t are part of the foreground block $\langle \cdot \rangle_t$. The extended graph operations are given in Figure 15.

A configuration graph G mirrors the structure of the thread pool μ ; it is a mapping from thread symbols to thread graphs:

$$G = a_1 \hookrightarrow g_1 \uplus \dots \uplus a_n \hookrightarrow g_n$$

The vertices, edges and foreground blocks of a configuration graph are the union of the vertices, edges and foreground blocks of the component thread graphs. If $a \hookrightarrow g_a \in G$, an edge (a, u, δ) may be viewed as an edge from the sink of g_a to u . If $g_a = \emptyset$, this edge is ignored. The metrics such as work, span and foreground width extend in the natural way to configuration graphs.

The judgment $\mu_l; \mu_g \Downarrow_{\mathcal{C}}^{\Delta} \{G\}$ generates a portion of a configuration graph from the threads in a partial thread pool μ_l by generating a thread graph for each thread and composing any non-empty graphs that result. As above, the whole thread pool μ_g is included so that threads may refer to other threads which are not currently under attention, but these threads are not included in G . If a thread is delayed with delay $\delta > 0$, its cost graph is composed serially after a fresh auxiliary vertex using an edge of weight δ .

The extended cost semantics allows us to assign costs (work, span, etc.) to programs, as represented by thread pools. The work and span of a thread pool that is in the middle of execution can be thought of as the *remaining* work and span of the program. The work of a thread pool μ under Δ is written $W(\mu, \Delta)$ and is defined as the maximum work over all dags that can be generated from μ :

$$W(\mu, \Delta) = \max\{W(G) \mid \mu; \mu \Downarrow_{\mathcal{C}}^{\Delta} \{G\}\}$$

We take the maximum since the cost semantics is nondeterministic. The definitions of $S(\mu, \Delta)$, $W^\circ(\mu, \Delta)$ and $S^\circ(\mu, \Delta)$ are similar.

For space reasons, we omit the proofs of two results. The first is the extension of Lemma 1 to handle programs that have begun to execute. The second is that the operational semantics and cost semantics agree on values produced by an expression. The main complication in showing the correspondence of the cost and operational semantics is that the value $\text{thread}[v](u)$ is produced by the cost semantics but not the operational semantics. We therefore show that the cost semantics and the operational semantics are equivalent up to a relation which relates the two forms of thread handle.

4.3 Cost Bounds for Prompt Scheduling Principle

The main result of this section is showing that the cost bounds predicted by the cost semantics can be realized by the operational semantics in that, given a prompt schedule, a λ^{ip} program can be evaluated using the operational semantics in the number of steps and response time predicted by the prompt scheduling theorem (Theorem 1).

The key step in showing the bound on the computation time is showing that a global transition step decreases the total work by P or the total span by 1. The intuition behind this proof is the same as that of Theorem 1: the scheduler will either execute P (foreground) instructions or execute all ready (foreground) instructions. The proof of this lemma makes heavy use of a technical lemma which shows that if $e \mid \mu \mapsto_a^\Delta (\delta, e') \mid \mu \uplus \mu'$ and $e; \mu \Downarrow^\Delta v; g$ and $e'; \mu \uplus \mu' \Downarrow^\Delta v'; g'$, then g' is the same as g with its source vertex removed. That is, the local transition decreases the work and span of the thread's dag by at least 1.

Lemma 2. *Fix Δ and suppose that $\cdot \vdash \mu : \Sigma$ and that e is well-joined for all $a \hookrightarrow (\delta, e) \in \mu$. If $r; \mu \mapsto_{g10} r'; \mu'$ using a prompt scheduling policy, then*

1. $W(\mu', \Delta) \leq W(\mu, \Delta)$
2. $S(\mu', \Delta) \leq S(\mu, \Delta)$
3. $W(\mu, \Delta) - W(\mu', \Delta) \geq P$ or $S(\mu, \Delta) - S(\mu', \Delta) \geq 1$
4. $W^\circ(\mu', \Delta) \leq W^\circ(\mu, \Delta)$
5. $S^\circ(\mu', \Delta) \leq S^\circ(\mu, \Delta)$
6. $W^\circ(\mu, \Delta) - W^\circ(\mu', \Delta) \geq P$ or $S^\circ(\mu, \Delta) - S^\circ(\mu', \Delta) \geq r' - r$.

Proof. See the companion technical report [46]. \square

The proof of the response time and computation time bounds is then straightforward.

Theorem 3. *Fix Δ and let e be such that $\cdot \vdash e : \tau @ \mathbb{B}$. Suppose $e; \emptyset \Downarrow^\Delta v; g$ If $0; a \hookrightarrow (0, e) \mapsto_{g10}^T r; \mu$ using a prompt scheduling policy and μ final, then $T \leq \frac{W(g)}{P} + S(g)$ and $r \leq D(g) \frac{W^\circ(g)}{P} + S^\circ(g)$.*

Proof. Let $\mu_0 = a \hookrightarrow (0, e)$ and $\mu_T = \mu$ and $r_0 = 0$ and $r_T = r$. We have a sequence $0; \mu_0 \mapsto_{g10} r_1; \mu_1 \mapsto_{g10} \dots \mapsto_{g10} r_T; \mu_T$.

For each i , let $W_i = W(\mu_i, \Delta)$ (and similar for S_i , W_i° and S_i°). Note that $W_0 = W(g)$ (and similar for S , W° , S°) and that $W_T = S_T = W_T^\circ = S_T^\circ = 0$.

By Lemma 2 and preservation of well-joinedness,

$$\frac{W_0}{P} + S_0 \geq 1 + \frac{W_1}{P} + S_1 \geq \dots \geq 1 + \frac{W_T}{P} + S_T = 1$$

This immediately gives $\frac{W_0}{P} + S_0 \geq T$.

Let $D = D(g)$. For each i , consider the quantity $D \frac{W_i^\circ}{P} + S_i^\circ + r_i$. Note that for $i = 0$, $D \frac{W_0^\circ}{P} + S_0^\circ + r_0 = D \frac{W^\circ(g)}{P} + S^\circ(g)$ and for $i = T$, $D \frac{W_T^\circ}{P} + S_T^\circ + r_T = r$. When $r_i; \mu_i \mapsto_{g10} r_{i+1}; \mu_{i+1}$, by Lemma 2, either

1. $W_i^\circ - W_{i+1}^\circ \geq P$ and $r_{i+1} - r_i = |RFB(\mu_i)| \leq D$ (the last inequality is by definition of D) or
2. $S_i^\circ - S_{i+1}^\circ \geq |RFB(\mu_i)|$ and $r_{i+1} - r_i = |RFB(\mu_i)|$

In both cases, the quantity above decreases or remains the same, so $r \leq D(g) \frac{W_0^\circ}{P} + S_0^\circ$. \square

5. Implementation and Evaluation

The operational semantics (Section 4) specifies an implementation at the level of threads and scheduling decisions. To realize the semantics in practice, we must implement the global scheduling step by giving a prompt scheduling algorithm. In order to approximate the operational semantics, which reschedules at each step, it is necessary to perform some preemption, using periodic interrupts, so that low-priority threads can be switched out for high-priority threads.

Next, prompt scheduling requires that, whenever the scheduler runs, it maps high-priority threads onto the available processors, followed by low-priority threads if any processors remain. A naive implementation could use a global priority queue, but this would not scale beyond just a few processors due to the cost of synchronization at the queue. A realistic implementation therefore would have to distribute the queues. There are many ways to achieve this. In this paper, we build on a recently proposed variant of the work-stealing algorithm [2]. In our algorithm, each processor has a private priority queue and a public communication cell, a *mailbox*, to which other processors can send threads. At periodic intervals, each processor attempts to send, or *deal*, a thread to a random processor, in priority order, by atomically writing into the target processor's mailbox. Each processor then checks its own queue and mailbox and begins working on the highest-priority task available. Generalizations of work stealing to support priorities have been considered before [34] but these algorithms are not preemptive.

5.1 Implementation

We implemented the basic primitives of our formal language λ^{ip} as a Standard ML library, and implemented the preemptive

priority-based work stealing algorithm described above by building on an existing parallel extension [58, 59] of the MLton [44] compiler for Standard ML. We have not extended SML’s type system to implement λ^{ip} ’s temporal type system because this is less essential for the performance analysis.

5.2 Experimental Setup

The experiments were performed on a 48-core machine with 125GB of memory and 2.1 GHz AMD CPUs running Ubuntu 16.04. To account for inherent noise in the data, we performed each run between 10 and 20 times and each data point represents the average over the runs. Through empirical analysis, we found that interrupt intervals in the range of 1-25 milliseconds lead to the best throughput and responsiveness. For the results reported in this paper, we use a 5ms interval.

Measuring Responsiveness. Empirical analysis of interactive programs can be challenging because it requires isolating the completion time of potentially small pieces of computation (such as an interaction with a user) within an application. For example, prior work proposed operating system modifications [21]. We use a simpler approach. In our experiments, a *driver* program, written in C, reads a sequence of interactive events (e.g. mouse clicks, key presses) from a trace file. It simulates these events, records the response of the program to the event and measures the response time.

5.3 Quantitative Benchmarks

Fibonacci-Terminal. This benchmark performs a parallel Fibonacci computation, specifically `fib(45)` (to stand in for an intensive parallel computation), and simultaneously performs user interaction via a terminal. The user interaction consists of a loop that repeatedly reads a name from standard input, and immediately greets the user by name. To ensure responsiveness, the benchmark designates the terminal computation as high priority and the Fibonacci computation as low priority. The benchmark terminates once the `fib(45)` computation completes.

To assess the responsiveness of the Fibonacci-terminal benchmark, we run it while varying the number of processors between 1 and 30 and the rate of interaction between 1 and 50 interactions per second⁵. In the experiments, the driver program sends a name on standard input at uniform intervals to match the desired number of interactions per second. It then waits for the response from the program. The time between the input and the response is the response time. We measure the response time for each input and take the average.

The left plot in Figure 16 shows the speedup (with respect to the sequential version of Fibonacci) of the Fibonacci computation as a function of the number of processors under varying interaction rates. For comparison, the figure also shows (in blue squares) the speedup of a standard work stealing scheduler running a Fibonacci computation only (with

⁵ Due to a technical limitation of the thread-pinning library used by the runtime, we were unable to use all of the system’s cores.

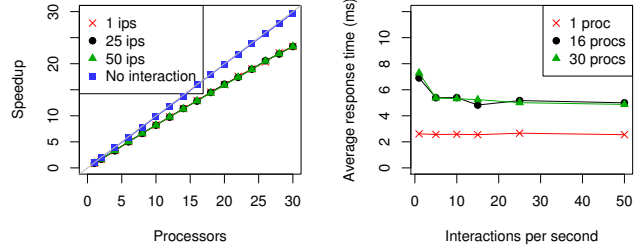


Figure 16. Speedup (l) and response time (r) for the Fibonacci-terminal benchmark.

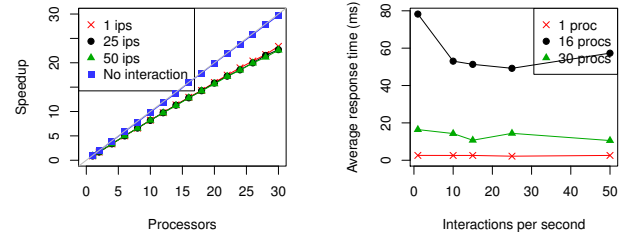


Figure 17. Speedup (l) and response time (r) for the Fibonacci-network benchmark.

no interaction). The results show that interaction decreases the speedup, but not significantly. This is consistent with our bounds because interaction, which is high-priority, takes precedence over the low-priority Fibonacci computation. The right plot in Figure 16 shows the average response time as a function of the number of interactions per second. The average response time remains relatively flat even as the interaction rate increases, which is expected because each interaction involves little work (just echoing the input name). We furthermore see that increasing the number of processors causes an increase in the response time up to a point. This seems counterintuitive but is likely caused by migrations of high-priority computations to other processors via a deal, which can increase response time compared to the local handling of the same interaction. Overall, average response time remains very good, staying well under 10 milliseconds.

Fibonacci-Network. Our next benchmark has the same structure as the Fibonacci-terminal but involves more complex interaction. The benchmark opens a socket and listens for incoming connections. When a connection is received, it starts an interactive channel, implemented as a new high-priority thread. The interaction on each channel proceeds as in Terminal echo above, until the program terminates or the client disconnects. Because there can be many channels, each of which is handled by a thread, active at the same time, this benchmark tests the case where there are many interactive computations, all of which demand responsiveness.

In this benchmark, the driver program opens a number of network connections, and sends a line over each at one-second intervals, staggered so that the messages arrive at uniform intervals. The number of network connections opened

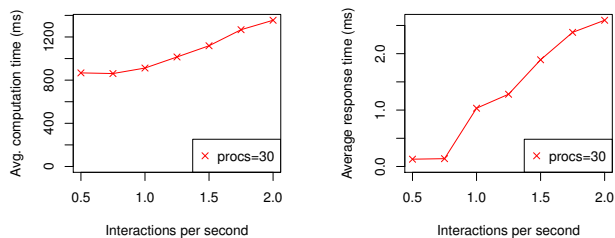


Figure 18. The effect of interaction rate for the Fibonacci server on: (l) Computation time (r) Response time.

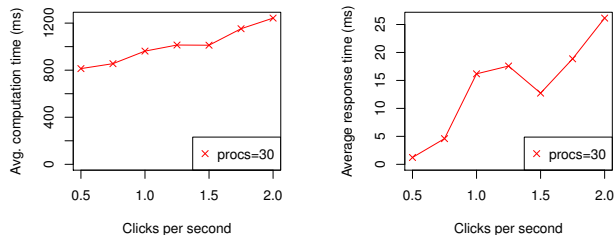


Figure 19. The effect of interaction rate for the convex hull server on: (l) Computation time (r) Response time.

is the desired number of interactions per second. Figure 17 shows the results. We see again that the Fibonacci computation scales well with respect to the sequential baseline with varying levels of interaction. As above, the one-processor case shows the best responsiveness and the average response times are good, under 100 milliseconds.

Fibonacci Server. In the above benchmarks, the interactive and computational parts of the benchmark did not interact, apart from the fact that they run together. Our next benchmark, the “Fibonacci server”, simulates an application that receives queries, each of which requires performing some compute-intensive task. An interactive, high-priority loop waits for the user to enter a number on the console. When a number n is entered, the loop starts the computation of the n^{th} Fibonacci number in a separate low-priority thread that also prints the result on the console; the loop continues to listen to further inputs immediately after starting the Fibonacci computation.

To assess this benchmark, our driver program runs the benchmark with a trace that inputs the numbers 41 to 45 in increasing order. The driver calculates the response time as the time between the input and the next prompt, and the computation time as the time taken to compute each Fibonacci number. The rate of interaction varies from 0.5 to 2.0 inputs per second. Figure 18 shows the results. As expected, both numbers begin to increase as the interaction becomes frequent enough that the Fibonacci computations overlap. The computations still complete in a timely matter, and the program remains responsive, with average response times not exceeding several milliseconds.

Interactive Convex Hull. Our interactive convex hull benchmark maintains the convex hull of a set of 2D points,

as a user inserts new points by clicking on the screen. A high-priority loop polls the mouse and, every time the user adds a point, starts a low-priority thread that computes and draws the new convex hull. In our experiments, the driver program simulates five clicks at random points on the screen at regular intervals and calculates the response time for each click as the time between the click and drawing of the point, and the computation time as the time to compute each hull. So that the hull computations are not trivial, each one includes 1,000,000 random points that are chosen at initialization. Figure 19 shows the results. As with the Fibonacci server, the computation times and response times increase with the interaction rate; this is expected because increased interaction rate causes multiple convex hull computations to overlap with each other and with the interaction. The program still remains responsive to clicks, with response times under 30ms.

5.4 Qualitative Benchmarks

In addition to the relatively simple benchmarks above, we considered more sophisticated benchmarks. These benchmarks are more difficult to evaluate quantitatively, but we assess their performance qualitatively by their usability.

Web Server. A high-priority loop listens for connections and starts a new high-priority thread for each one. HTTP requests are logged, and a low-priority thread periodically performs analytics on the log. We simulate a large analytics computation by computing a large Fibonacci number in parallel. As expected, we observe that the background computations do not interfere with the handling of HTTP requests.

Photo Viewer. Our photo viewer benchmark allows the user to navigate through a folder of JPEG images, either by scrolling or jumping to an image. To ensure smooth scrolling, the user interaction is high priority, and the viewer decodes the next several images in the background so they will be ready when requested. If the user selects an image that has not yet been decoded, it is decoded in the foreground and displayed. Our experience shows that the viewer is responsive, indicating that the decoding is proceeding quickly enough to be effective, and that the background decoding processes do not hamper the high-priority interaction.

Music Server. A streaming music server listens for network connections and spawns a thread for each new client. The client requests a music file from the server, which the server streams over the connection until the end of the file is reached. Some clients (perhaps those paying for a higher level of service) are designated high-priority, and are handled by high-priority threads; the remaining clients receive low priority. We tested the server with a relatively small number of clients (up to 10, both low and high priority), and in our experience, it maintains a high quality of service for all clients.

6. Related Work

We discussed the most closely related work in the main body of the paper. Here we take a broader perspective and briefly describe more remotely related work.

Parallel Computing. Much work has been done on parallel computing with dynamically scheduled, fine-grained and cooperative threads since the 1970s [5, 10, 17, 18, 25, 26, 30, 35–37, 39, 41, 42, 51]. Nearly all of this work focuses on maximizing throughput in compute-intensive applications and relies on cooperative threading. This paper shows that the language abstractions, dag-based cost models [13, 28, 38] and cost semantics [8, 9, 29], can be extended to include competitive threading, where threads are scheduled preemptively.

Type Systems for Staged Computation. The type system of λ^{ip} is based on that of Davies [19] for binding time analysis, which is derived from linear temporal logic. This work influenced much followup work on metaprogramming and staged computation [40, 48, 50, 61]. These systems allow a computation at a stage to create and manipulate, but not eliminate, a computation in a later stage. For example, a stage 1 computation can create a stage 2 computation as a “black box” but cannot inspect that computation. We use a two-stage variant of the \circ modality of Davies [19], similar to that of Feltman et al. [23], which inspires some of our notation. One important difference between stages and the priorities of our work is that, in our work, computations belonging to different stages (priorities) can be evaluated concurrently, whereas in staged computations, evaluation proceeds monotonically in stage order.

Cost Semantics. The idea of using a cost semantics to reason about efficiency of programs goes back to the early 1990s [52, 54] and has since been applied in a number of contexts [8, 9, 43, 54, 55, 59]. Our approach builds directly on the work of Blleloch and Greiner [9] and Spoonhower et al. [59], who use computation graphs represented as dags (directed acyclic graphs) to reason about time and space in functional parallel programs. These cost models, however, consider cooperatively threaded parallelism only.

Scheduling. Our prompt-scheduling results generalize Brent’s classic result for scheduling parallel computations [15]. Since Brent’s result, much work has been done on scheduling. Ullman [62], Brent [15], and Eager et al. [20] established the hardness of optimal scheduling and the greedy scheduling principle. These early results have led to many more algorithms [1–3, 13, 16, 22, 26, 30, 49]. More recent papers showed that priority-based schedulers can improve performance in practice [34, 63, 64]. Our weighted-dag model builds on the model of Muller and Acar [45], who developed an algorithm for scheduling blocking parallel programs to hide latency, but did not consider responsiveness.

Scheduling is also studied extensively in the operating systems community (a book by Silberschatz et al. [56] presents

a comprehensive overview). There has been significant interest in making operating systems work well on multicore machines [6, 14]. The focus, however, has been on reducing contention within the OS and, as in the high-performance computing community, distributing resources to jobs so that they can run effectively. Scheduling within a job, which is our main concern, is less central to systems research.

There has been a great deal of work on scheduling for responsiveness in queuing theory (Harchol-Balter [31] presents a comprehensive overview). This line of work assumes a continuous stream of independent jobs arriving for processing according to some stochastic process. Such arrival assumptions do not quite fit the parallel computing model, where work is created by a program. In queuing theory, each job is generally processed by a single processor (or “server”) that decides at every point in time which of the current jobs to run. This work, however, typically assumes jobs to be sequential.

Scheduling is also an important concern in real-time computing. Most of this work considers highly structured (usually synchronous) sequential computations. Saifullah et al. [53] consider scheduling a set of real-time tasks where each task is a parallel computation represented by a parallel dag. Their algorithm infers for each vertex in the dag a deadline and schedules the vertices according to their deadlines. Their work assumes that the tasks are independent and are known in advance, as is the dag structure.

7. Conclusion

This paper takes a step toward uniting cooperative and competitive threading. To this end, we consider a programming language with fork-join parallelism, interaction, and priorities, and extend the classic cost models for cooperative threading based on cost graphs and cost semantics to bound both run-time and responsiveness. Our implementation and experiments suggest that the approach can be made practical. We leave a number of questions to future work, including the extension of our techniques to multiple priorities (instead of the two priorities we consider), the development of an efficient scheduling algorithm that implements the prompt-scheduling principle, and a more detailed evaluation.

Acknowledgments

This research is partially supported by grants from the National Science Foundation (CCF-1320563, CCF-1408940, CCF-1629444) and European Research Council (grant ERC-2012-StG-308246), and by a gift from Microsoft Research. We thank Tim Harris and Ziv Scully for their feedback.

References

- [1] U. A. Acar, G. E. Blleloch, and R. D. Blumofe. The data locality of work stealing. *Theory of Computing Systems (TOCS)*, 35(3):321–347, 2002.

- [2] U. A. Acar, A. Charguéraud, and M. Rainey. Scheduling parallel programs by work stealing with private dequeues. In *PPoPP '13*, 2013.
- [3] U. A. Acar, A. Charguéraud, and M. Rainey. Oracle-guided scheduling for controlling granularity in implicitly parallel languages. *Journal of Functional Programming (JFP)*, 26:e23, 2016.
- [4] N. S. Arora, R. D. Blumofe, and C. G. Plaxton. Thread scheduling for multiprogrammed multiprocessors. *Theory of Computing Systems*, 34(2):115–144, 2001.
- [5] Arvind and K. P. Gostelow. The Id report: An asynchronous language and computing machine. Technical Report TR-114, Department of Information and Computer Science, University of California, Irvine, Sept. 1978.
- [6] A. Baumann, P. Barham, P.-E. Dagand, T. Harris, R. Isaacs, S. Peter, T. Roscoe, A. Schüpbach, and A. Singhanian. The multikernel: A new OS architecture for scalable multicore systems. In *Proceedings of the ACM SIGOPS 22nd Symposium on Operating Systems Principles, SOSP '09*, pages 29–44, 2009.
- [7] G. Blake, R. G. Dreslinski, T. Mudge, and K. Flautner. Evolution of thread-level parallelism in desktop applications. In *Proceedings of the 37th Annual International Symposium on Computer Architecture, ISCA '10*, pages 302–313, 2010.
- [8] G. Blleloch and J. Greiner. Parallelism in sequential functional languages. In *Proceedings of the 7th International Conference on Functional Programming Languages and Computer Architecture, FPCA '95*, pages 226–237. ACM, 1995.
- [9] G. E. Blleloch and J. 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.
- [10] G. E. Blleloch, J. C. Hardwick, J. Sipelstein, M. Zagha, and S. Chatterjee. Implementation of a portable nested data-parallel language. *J. Parallel Distrib. Comput.*, 21(1):4–14, 1994.
- [11] G. E. Blleloch, P. B. Gibbons, and Y. Matias. Provably efficient scheduling for languages with fine-grained parallelism. *J. ACM*, 46:281–321, Mar. 1999.
- [12] R. D. Blumofe and C. E. Leiserson. Space-efficient scheduling of multithreaded computations. *SIAM Journal on Computing*, 27(1):202–229, 1998.
- [13] R. D. Blumofe and C. E. Leiserson. Scheduling multithreaded computations by work stealing. *J. ACM*, 46:720–748, Sept. 1999.
- [14] S. Boyd-Wickizer, H. Chen, R. Chen, Y. Mao, F. Kaashoek, R. Morris, A. Pesterev, L. Stein, M. Wu, Y. Dai, Y. Zhang, and Z. Zhang. Core: An operating system for many cores. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation*, pages 43–57, 2008.
- [15] R. P. Brent. The parallel evaluation of general arithmetic expressions. *J. ACM*, 21(2):201–206, 1974.
- [16] F. W. Burton and M. R. Sleep. Executing functional programs on a virtual tree of processors. In *Functional Programming Languages and Computer Architecture (FPCA '81)*, pages 187–194. ACM Press, Oct. 1981.
- [17] M. M. T. Chakravarty, R. Leshchinskiy, S. L. Peyton Jones, G. Keller, and S. Marlow. Data parallel Haskell: a status report. In *Proceedings of the POPL 2007 Workshop on Declarative Aspects of Multicore Programming, DAMP 2007, Nice, France, January 16, 2007*, pages 10–18, 2007.
- [18] P. Charles, C. Grothoff, V. Saraswat, C. Donawa, A. Kielstra, K. Ebcioglu, C. von Praun, and V. Sarkar. X10: an object-oriented approach to non-uniform cluster computing. In *Proceedings of the 20th annual ACM SIGPLAN conference on Object-oriented programming, systems, languages, and applications, OOPSLA '05*, pages 519–538. ACM, 2005.
- [19] R. Davies. A temporal-logic approach to binding-time analysis. In *LICS*, pages 184–195, 1996.
- [20] D. L. Eager, J. Zahorjan, and E. D. Lazowska. Speedup versus efficiency in parallel systems. *IEEE Transactions on Computing*, 38(3):408–423, 1989.
- [21] Y. Endo, Z. Wang, J. B. Chen, and M. Seltzer. Using latency to evaluate interactive system performance. In *Proceedings of the Second USENIX Symposium on Operating Systems Design and Implementation, OSDI '96*, pages 185–199, New York, NY, USA, 1996. ACM.
- [22] D. G. Feitelson, L. Rudolph, and U. Schwiegelshohn. Parallel job scheduling - A status report. In *Job Scheduling Strategies for Parallel Processing (JSSPP), 10th International Workshop*, pages 1–16, 2004.
- [23] N. Feltman, C. Angiuli, U. A. Acar, and K. Fatahalian. Automatically splitting a two-stage lambda calculus. In *Proceedings of the 25 European Symposium on Programming, ESOP*, pages 255–281, 2016.
- [24] K. Flautner, R. Uhlig, S. Reinhardt, and T. Mudge. Thread-level parallelism and interactive performance of desktop applications. In *Proceedings of the Ninth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS IX*, pages 129–138, 2000.
- [25] M. Fluett, M. Rainey, J. Reppy, and A. Shaw. Implicitly threaded parallelism in Manticore. *Journal of Functional Programming*, 20(5-6):1–40, 2011.
- [26] M. Frigo, C. E. Leiserson, and K. H. Randall. The implementation of the Cilk-5 multithreaded language. In *PLDI*, pages 212–223, 1998.
- [27] C. Gao, A. Gutierrez, R. G. Dreslinski, T. Mudge, K. Flautner, and G. Blake. A study of thread level parallelism on mobile devices. In *Performance Analysis of Systems and Software (ISPASS), 2014 IEEE International Symposium on*, pages 126–127, March 2014.
- [28] R. L. Graham. Bounds on multiprocessing timing anomalies. *SIAM Journal on Applied Mathematics*, 17(2):416–429, 1969.
- [29] J. Greiner and G. E. Blleloch. A provably time-efficient parallel implementation of full speculation. *ACM Transactions on Programming Languages and Systems*, 21(2):240–285, Mar. 1999.
- [30] R. H. Halstead. Multilisp: a language for concurrent symbolic computation. *ACM Transactions on Programming Languages and Systems*, 7:501–538, 1985.
- [31] M. Harchol-Balter. *Performance Modeling and Design of Computer Systems: Queuing Theory in Action*. Cambridge

University Press, 2013.

- [32] R. Harper. *Practical Foundations for Programming Languages*. Cambridge University Press, New York, NY, USA, 2012.
- [33] C. Hauser, C. Jacobi, M. Theimer, B. Welch, and M. Weiser. Using threads in interactive systems: A case study. In *Proceedings of the Fourteenth ACM Symposium on Operating Systems Principles*, SOSOP '93, pages 94–105, New York, NY, USA, 1993. ACM.
- [34] S. Imam and V. Sarkar. Load balancing prioritized tasks via work-stealing. In *Euro-Par 2015: Parallel Processing - 21st International Conference on Parallel and Distributed Computing*, pages 222–234, 2015.
- [35] S. M. Imam and V. Sarkar. Habanero-Java library: a Java 8 framework for multicore programming. In *2014 International Conference on Principles and Practices of Programming on the Java Platform Virtual Machines, Languages and Tools, PPPJ '14*, pages 75–86, 2014.
- [36] Intel. Intel threading building blocks, 2011. <https://www.threadingbuildingblocks.org/>.
- [37] S. Jagannathan, A. Navabi, K. Sivaramakrishnan, and L. Ziarek. The design rationale for Multi-MLton. In *ML '10: Proceedings of the ACM SIGPLAN Workshop on ML*. ACM, 2010.
- [38] J. Jaja. *An introduction to parallel algorithms*. Addison Wesley Longman Publishing Company, 1992.
- [39] G. Keller, M. M. Chakravarty, R. Leshchinskiy, S. Peyton Jones, and B. Lippmeier. Regular, shape-polymorphic, parallel arrays in Haskell. In *Proceedings of the 15th ACM SIGPLAN international conference on Functional programming*, ICFP '10, pages 261–272, 2010.
- [40] T. B. Knoblock and E. Ruf. Data specialization. In *Proceedings of the ACM SIGPLAN 1996 Conference on Programming Language Design and Implementation*, PLDI '96, pages 215–225, 1996.
- [41] D. Lea. A Java fork/join framework. In *Proceedings of the ACM 2000 conference on Java Grande, JAVA '00*, pages 36–43, 2000.
- [42] D. Leijen, W. Schulte, and S. Burckhardt. The design of a task parallel library. In *Proceedings of the 24th ACM SIGPLAN conference on Object Oriented Programming Systems Languages and Applications*, OOPSLA '09, pages 227–242, 2009.
- [43] R. Ley-Wild, U. A. Acar, and M. Fluet. A cost semantics for self-adjusting computation. In *Proceedings of the 26th Annual ACM Symposium on Principles of Programming Languages*, POPL '09, 2009.
- [44] MLton. MLton web site. <http://www.mlton.org>.
- [45] S. K. Muller and U. A. Acar. Latency-hiding work stealing: Scheduling interacting parallel computations with work stealing. In *Proceedings of the 28th ACM Symposium on Parallelism in Algorithms and Architectures*, SPAA '16, pages 71–82, 2016.
- [46] S. K. Muller, U. A. Acar, and R. Harper. Responsive parallel computation: Bridging competitive and cooperative threading. Technical Report TBD, Carnegie Mellon University School of Computer Science, Apr. 2017.
- [47] T. Murphy, VII, K. Crary, R. Harper, and F. Pfenning. A symmetric modal lambda calculus for distributed computing. In *Proceedings of the 19th IEEE Symposium on Logic in Computer Science (LICS)*, pages 286–295. IEEE Press, 2004.
- [48] A. Nanevski and F. Pfenning. Staged computation with names and necessity. *J. Funct. Program.*, 15(5):893–939, 2005.
- [49] G. J. Narlikar and G. E. Blelloch. Space-efficient scheduling of nested parallelism. *ACM Transactions on Programming Languages and Systems*, 21, 1999.
- [50] F. Pfenning and R. Davies. A judgmental reconstruction of modal logic. *Mathematical Structures in Computer Science*, 11:511–540, 2001.
- [51] R. Raghunathan, S. K. Muller, U. A. Acar, and G. Blelloch. Hierarchical memory management for parallel programs. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming*, ICFP 2016, pages 392–406, New York, NY, USA, 2016. ACM.
- [52] M. Rosendahl. Automatic complexity analysis. In *FPCA '89: Functional Programming Languages and Computer Architecture*, pages 144–156. ACM, 1989.
- [53] A. Saifullah, D. Ferry, J. Li, K. Agrawal, C. Lu, and C. D. Gill. Parallel real-time scheduling of DAGs. *IEEE Trans. Parallel Distrib. Syst.*, 25(12):3242–3252, 2014.
- [54] D. Sands. Complexity analysis for a lazy higher-order language. In *ESOP '90: Proceedings of the 3rd European Symposium on Programming*, pages 361–376, London, UK, 1990. Springer-Verlag.
- [55] P. M. Sansom and S. L. Peyton Jones. Time and space profiling for non-strict, higher-order functional languages. In *Principles of Programming Languages*, pages 355–366, 1995.
- [56] A. Silberschatz, P. B. Galvin, and G. Gagne. *Operating system concepts (7. ed.)*. Wiley, 2005.
- [57] D. C. Smith, C. Irby, R. Kimball, B. Verplank, and E. Harslem. Designing the star user interface. *BYTE Magazine*, 7(4):242–282, 1982.
- [58] D. Spoonhower. *Scheduling Deterministic Parallel Programs*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, USA, 2009.
- [59] D. Spoonhower, G. E. Blelloch, R. Harper, and P. B. Gibbons. Space profiling for parallel functional programs. In *International Conference on Functional Programming*, 2008.
- [60] D. C. Swinehart, P. T. Zellweger, R. J. Beach, and R. B. Hagmann. A structural view of the cedar programming environment. *ACM Trans. Program. Lang. Syst.*, 8(4):419–490, Aug. 1986.
- [61] W. Taha and T. Sheard. MetaML and multi-stage programming with explicit annotations. *Theoretical Computer Science*, 248(1):211 – 242, 2000.
- [62] J. Ullman. NP-complete scheduling problems. *Journal of Computer and System Sciences*, 10(3):384 – 393, 1975.
- [63] M. Wimmer, D. Cederman, J. L. Träff, and P. Tsigas. Work-stealing with configurable scheduling strategies. In *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPoPP '13, pages 315–316, 2013.

[64] M. Wimmer, F. Versaci, J. L. Träff, D. Cederman, and P. Tsigas. Data structures for task-based priority scheduling. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and*

Practice of Parallel Programming, PPOPP '14, pages 379–380, 2014.