

# Task Parallel Assembly Language for Uncompromising Parallelism

Mike Rainey  
Carnegie Mellon University  
Pittsburgh, PA, USA  
me@mike-rainey.site

Ryan R. Newton  
Facebook  
New York, NY, USA  
Newton@fb.com

Kyle Hale  
Illinois Institute of Technology  
Chicago, IL, USA  
khale@cs.iit.edu

Nikos Hardavellas  
Northwestern University  
Chicago, IL, USA  
nikos@northwestern.edu

Simone Campanoni  
Northwestern University  
Chicago, IL, USA  
simonec@northwestern.edu

Peter Dinda  
Northwestern University  
Chicago, IL, USA  
pdinda@northwestern.edu

Umut A. Acar  
Carnegie Mellon University  
Pittsburgh, PA, USA  
umut@cs.cmu.edu

## Abstract

Achieving parallel performance and scalability involves making compromises between parallel and sequential computation. If not contained, the overheads of parallelism can easily outweigh its benefits, sometimes by orders of magnitude. Today, we expect programmers to implement this compromise by optimizing their code manually. This process is labor intensive, requires deep expertise, and reduces code quality. Recent work on *heartbeat scheduling* shows a promising approach that manifests the potentially vast amounts of available, latent parallelism, at a regular rate, based on even beats in time. The idea is to amortize the overheads of parallelism over the useful work performed between the beats. Heartbeat scheduling is promising in theory, but the reality is complicated: it has no known practical implementation.

In this paper, we propose a practical approach to heartbeat scheduling that involves equipping the assembly language with a small set of primitives. These primitives leverage existing kernel and hardware support for interrupts to allow parallelism to remain latent, until a heartbeat, when it can be manifested with low cost. Our Task Parallel Assembly Language (TPAL) is a compact, RISC-like assembly language.

We specify TPAL through an abstract machine and implement the abstract machine as compiler transformations for C/C++ code and a specialized run-time system. We present an evaluation on both the Linux and the Nautilus kernels, considering a range of heartbeat interrupt mechanisms. The evaluation shows that TPAL can dramatically reduce the overheads of parallelism without compromising scalability.

**CCS Concepts:** • Software and its engineering → Parallel programming languages.

**Keywords:** parallel programming languages, granularity control

## ACM Reference Format:

Mike Rainey, Ryan R. Newton, Kyle Hale, Nikos Hardavellas, Simone Campanoni, Peter Dinda, and Umut A. Acar. 2021. Task Parallel Assembly Language for Uncompromising Parallelism. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI '21)*, June 20–25, 2021, Virtual, Canada. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3453483.3460969>

## 1 Introduction

A classic problem in parallel computing is to take a high-level parallel program, written in nested-parallel style, with fork-join constructs, and derive from it an executable that is efficient on real machines. Traditionally, solutions involve optimizing the program to control the amount of parallelism exposed, thereby limiting the overheads of task creation and scheduling [3, 7, 30, 60]. Left unchecked, task overheads can reach two orders of magnitude or more, effectively wiping out the benefits of parallelism. But when task overheads are addressed, the associated optimizations involve changing the code so that the program switches from parallel to sequential code, typically at “small” problem sizes. This is a

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org). PLDI '21, June 20–25, 2021, Virtual, Canada

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8391-2/21/06...\$15.00

<https://doi.org/10.1145/3453483.3460969>

process called *granularity control* [3, 7]. In addition to labor-intensive changes, granularity control requires tuning the program so that it switches from parallel to serial at appropriate points at run time. Such tuning may cause the program to lose its performance portability, because the process of tuning usually overfits to the idiosyncrasies of a particular machine. In particular, the notion of “small” depends on the machine architecture and software environment and varies significantly from one machine to another [60].

Motivated by the limitations of manual code optimizations and granularity control, recent work proposed an alternative approach that can, in principle, be completely automated. This approach, called heartbeat scheduling [5], provably controls the overheads of parallelism without requiring manual changes to the code. In heartbeat scheduling, a regular heartbeat event interrupts the program periodically to promote latent opportunities for parallelism into actual tasks that can be executed in parallel, e.g., by migrating across cores. Therefore, rather than making granularity decisions in the source code of the program, the program exposes the *maximum* amount of parallelism, and the heartbeat mechanism decides how and when to promote latent parallelism into actual parallelism. Intuitively speaking, because the heartbeat only fires at certain points in time, the cost of creating and managing parallelism can be amortized over the useful work done between heartbeats.

Although the idea of the heartbeat scheduling may appear quite simple, its realization is far from it. Perhaps the most important challenge is determining, at each heartbeat, a unit of work that must be reified as a task, so that it can be executed in parallel, just like any other task that is created by carefully hand-optimized codes. Prior work evaluates heartbeat scheduling by using a C++ interpreter that runs programs hand-rolled in a custom format, representing programs as abstract syntax trees. This structure allows the interpreter to create tasks by manipulating the abstract syntax tree at a heartbeat event. But runtime interpretive overhead is impractical, especially for high-performance parallel programming, and the interpreter provides only approximate timing of heartbeats.

In this paper, we propose an execution model for heartbeat scheduling, formalized by our Task Parallel Assembly Language (TPAL), and a runtime system that is backed by a practical heartbeat mechanism, which is implemented on top of hardware-based interrupts. The core of TPAL resembles a conventional, RISC assembly language. But unconventionally, TPAL features native support for task parallelism, consisting of a small collection of primitives and annotations on basic blocks. Because its task parallelism is native, TPAL can express parallel loops naturally and execute them efficiently. TPAL can achieve task parallelism as a *nearly zero-cost abstraction*, even with programs involving irregular, and nested parallel loops.

Because it is specified as an abstract machine, TPAL’s execution model is low-level and detailed enough to be thought of as a model for an implementation. To evaluate this execution model, we present an implementation of TPAL along with a runtime system supporting the operations needed for implementing heartbeat events and parallel evaluation. The runtime system is written in C++, and is optimized to achieve practical efficiency by using state-of-the-art scheduling techniques. To implement the heartbeat events, the runtime system relies on hardware interrupts, which can be controlled at greater precision and lower cost than by using software-based approximations. Because the heartbeat events are hardware-driven, the runtime system is reasonably decoupled from the specific implementations of heartbeat events, and is quite portable, allowing it to be used on different hardware and software platforms.

To evaluate the effectiveness of proposed techniques based on TPAL coupled with hardware interrupts, we consider a Linux-based system and an experimental kernel framework, called Nautilus [34]. Nautilus is specifically designed to support parallel runtime systems, and thus allows tight control over hardware resources and elides many of the overheads and abstraction layers that arise from supporting general-purpose workloads.

For our evaluation, we consider a number of parallel benchmarks, including those that exhibit high degrees of irregular parallelism. Our results indicate that TPAL, driven by hardware-based interrupts, achieves excellent work-efficiency, incurring small overheads for uniprocessor executions, without sacrificing parallelism, and scaling well as the number of cores increases. Compared to Cilk, a state-of-the-art task-parallel system, TPAL is comparable or better, even as it manages all parallelism *automatically*, driven by hardware interrupts. Our results from experimenting with Nautilus show that there is room for improvement in the efficiency of Linux-based, software signaling mechanisms.

This paper makes the following contributions:

- TPAL: a Task-Parallel Assembly Language for task-parallel programming;
- An execution model for TPAL as an abstract machine that controls creation of parallel tasks automatically using hardware driven interrupts;
- A portable runtime system for TPAL, written in C++, that can be used to run TPAL programs on modern multicore hardware;
- An evaluation both on Linux and Nautilus, an experimental kernel framework for parallel runtimes.

## 2 Task Parallel Assembly Language

For heartbeat scheduling, we need the ability to interrupt a running program and examine its runtime state efficiently. To see how, consider the following simple loop, written in C,

as a running example. It calculates the product of numbers passed in *a* and *b* (using addition), and leaves the result in *c*.

```
int a,b,c; int r = 0;
for (; a != 0; a--) { r += b; }
c = r;
```

To parallelize this loop with minimal cost, we ideally do not want to change *anything* about how the sequential version is compiled and executed. To this end, we propose to extend the assembly language with instructions that will allow executing the sequential code without any overheads, while, at the same time, making it possible to manifest latent parallelism. For example, the loop's index variable *a* is likely kept in a register rather than on the stack, and if we wish to take a heartbeat at runtime and *split* the remaining iterations (for parallelism), then we must find the current iteration in the appropriate register, rather than relying on finding that index in memory. This constraint bears resemblance to the one faced by a debugger, that is, in terms of interpreting runtime program state. But we cannot tolerate anything nearly as expensive as a debugger implementation (e.g., Linux *ptrace*).

Figure 2 shows the code for our running example in TPAL. In this code, if we assign the empty block annotations,  $\star_{exit} = \star_{loop} = \cdot$ , the first three blocks of the resulting code are identical to ones produced by a sequential C compiler. Apart from using jump statements, this code is similar to the C source code. Once this code executes, it is irrevocably sequential. But the unique ability of TPAL, among assembly languages<sup>1</sup>, is that we can add parallelism *without* changing the sequential behavior, and in fact without changing the code, except by adding annotations. We mark *promotable* program points, where it is possible to switch from the sequential loop to a parallel variant, and then join again afterwards. More precisely, we derive the parallel version by assigning  $\star_{exit} = \text{jtppt assoc-comm}; \{r \mapsto r2\}; \text{comb}$  and  $\star_{loop} = \text{prppt loop-try-promote}$ .

TPAL captures exactly the essential property of interruptibility, and provides the compiler with building blocks to construct semantically equivalent parallel and sequential variants of a function, with the ability to redirect between them at runtime. Because TPAL is a general-purpose, abstract assembly language, it can be targeted by a wide range of parallel source languages and can in turn be lowered to existing ISAs or low-level intermediate representations. (We target x86-64 in our evaluation, Section 4.)

## 2.1 Syntax and Execution Model

The syntax of TPAL is presented in Figure 1. It is based on a subset of the MIPS assembly language, using a familiar

$r \in \text{registers}, l \in \text{labels},$

$n \in \text{integer literals}, j \in \text{join-record identifiers}$

```
v ::= r | l | n | j
op ::= + | - | ...
l ::= r := v | r := op r, v | if-jump r, v
    | r := jralloc v | fork r, v
I ::= jump v | l; I | halt | join v
B ::= [★] I
★ ::= · | prppt l | jtppt jp; ΔR; l
jp ::= assoc | assoc-comm
ΔR ::= {r1 ↦ r'1, ..., rn ↦ r'n}
```

**Figure 1.** Grammar of TPAL. Highlighted syntax is specific to our parallel extensions, whereas the rest represents a conventional RISC instruction set.

```
1 prod: [·] // computes c = a * b
2   r := 0; jump loop
3 exit: [★exit]
4   c := r; halt
5 loop: [★loop]
6   if-jump a, exit;
7   r := r + b; a := a - 1;
8   jump loop
9 loop-try-promote: [·]
10  t := a < 2; if-jump t, loop;
11  jr := jralloc exit;
12  jump loop-promote
13 loop-par-try-promote: [·]
14  t := a < 2; if-jump t, loop-par;
15  jump loop-promote
16 loop-promote: [·]
17  m := a / 2; n := a % 2; a := m;
18  tr := r; r := 0;
19  fork jr, loop-par;
20  a := m + n; r := tr;
21  jump loop-par
22 loop-par: [prppt loop-par-try-promote]
23  if-jump a, exit-par;
24  r := r + b; a := a - 1;
25  jump loop-par
26 comb : [·]
27  r := r + r2; join jr
28 exit-par: [·]
29  join jr
```

**Figure 2.** The prod program in TPAL. Applying the empty block annotations,  $\star_{exit} = \star_{loop} = \cdot$ , yields a serial program. Applying  $\star_{exit} = \text{jtppt assoc-comm}; \{r \mapsto r2\}; \text{comb}$  and  $\star_{loop} = \text{prppt loop-try-promote}$ , yields a parallel program.

<sup>1</sup>The ability to elide parallelism annotations and have a semantically equivalent program is a goal shared by spawn/sync annotations in Cilk, or parallelism combinators in Haskell, but at the assembly level it is quite a different thing.

notation for instructions for readability. In TPAL, the execution of a program consists of a set of concurrent tasks, such that each task has its own private register file and call stack. Heap memory can be shared. How the tasks themselves are scheduled, that is, the order in which tasks run (up to dependency constraints) and on which core they run, is up to the *load-balancing algorithm*. TPAL is agnostic to load-balancing algorithm: it is compatible with e.g., variants of work stealing [6, 16, 19] or parallel depth first [46].

We start with the subset of the language that supports register-based memory, and later address the stack and heap. Our assembly language assumes a set of registers,  $r$ , labels  $l$ , integer literals  $n$ , and a special set of values that we call *join records*. A join record  $j$  is memory that is used by TPAL programs to synchronize multiple tasks at a join point. An operand  $v$  is either a register, a label, an integer literal, or a join record. A primitive operation  $op$  is one of a number of primitive operations that can be found on a conventional RISC machine, such as arithmetic operations for integers.

An instruction  $i$  is either a move-to-register operation, a primitive operation, a conditional jump, a join-record allocation, or a fork instruction. A join-record allocation instruction allocates (on the heap) and initializes a new join record. Its argument is a label that is to be the *continuation block* of the join point. A fork instruction spawns a new task, in a fashion resembling that of the UNIX `fork()` system call. It takes two arguments: first, a join record, and, second, a label from which the spawned task is to start executing. We call the spawned task the *child task* and the calling task the *parent task*. When it executes, the first action of the fork instruction is to register the dependency edge between the parent and child task in the join record. After the dependency is registered, the child task is added to the set of executing tasks, from which point it starts executing with a copy of the register file of its parent. The child task starts by executing the block at the label passed for the second argument of the fork instruction. After it issues the fork instruction, the parent task proceeds to issue its next instruction.

An instruction sequence  $I$  is a list-based representation of a sequence of assembly instructions: it is either a jump instruction, a sequencing operator (semicolon), a halt instruction, or a join instruction. A jump instruction is an unconditional jump operation. A semicolon operator specifies a sequential order between an instruction and an instruction sequence. A halt instruction terminates the whole machine. A join instruction initiates synchronization between a parent and one of its child tasks. Its first argument is a join record. When it executes this instruction, a task participates in a *join-resolution policy*. A join-resolution policy specifies the manner in which to combine the results held in the memories of the parent and child tasks. Upon completion of a join, that is, after all tasks registered in the join record issue their join instructions, the program then jumps to the label originally passed in the allocation of the join record.

A program is represented by a set of (labeled) code blocks. Each code block consists of an instruction sequence, along with an annotation. Such annotations, denoted by  $\star$ , are either an empty annotation, a promotion-ready program point, or a join-target program point. A *promotion-ready program point* is the entry point of a block for which there is a special behavior: when control targets the block, either control can flow, as usual, into the first instruction of the block, or control can flow, instead, to the label attached to the annotation. A *join-target program point* is the entry point of a code block that is assigned to be the continuation of the join point of parallel tasks. The join-target annotation specifies the join-resolution policy, i.e., the instruction sequence to be executed upon parent and child tasks meeting at their common join points. Its first component  $jp$  specifies whether the combining operation is only associative or both associative and commutative. Its second component  $\Delta R$  specifies the way in which joining tasks combine their register files into one register file, which is to be used by the combining block. Its third component specifies the label of the combining block.

## 2.2 Dynamics

We designed TPAL to make it natural and efficient to implement heartbeat scheduling: parallelism is introduced on a regular basis, in a two-stage process. The first stage involves the triggering of an interrupt and the second involves the manifestation of latent parallelism by the interrupt handler, which may fork a new task. The triggering of an interrupt is supported in TPAL by assigning each task a cycle counter, which increments every time a task issues an instruction. When the cycle counter of a task exceeds a certain threshold, that task is ready to trigger an interrupt. The threshold is a global parameter, written  $\heartsuit$ , and is determined by a one-time, per-machine tuning process, which is required by heartbeat scheduling [5]. The setting of the parameter is picked by the heartbeat tuner application to be just large enough to amortize the creation of a new task, but small enough to avoid pruning away useful amounts of parallelism.

**Adding Parallelism.** We return to our running example program, `prod`. By using the parallel block annotations, the program will ultimately allow the loop to be parallelized on demand. The serial-by-default structure in the program is its main strength: it is the reason we can achieve near zero-cost abstraction.

**Heartbeat interrupts.** When the heartbeat threshold is exceeded by a task, an interrupt is ready to be serviced the next time the program enters a promotion-ready program point. In our example, every time a task enters the loop block and its heartbeat threshold has passed, the task jumps to its handler block, namely `loop-try-promote`, instead of the first instruction of the loop block. The handler block checks on line 10 if there is any parallelism available in the remaining iterations of the loop. If there is no latent



parallelism, then our running task jumps back to the loop block, where it left off, but if there is, the task manifests the parallelism.

**Promotion.** To manifest the parallelism from the loop, our task performs a *promotion*, wherein a task creates a join point and forks a new, child task. Promotion begins in our running example on line 11, where our task allocates a join point. By passing to the join-allocation instruction the `exit` label, we are instructing our newly parallelized loop to terminate, with the final result, by jumping to the `exit` block. The handler next jumps to the loop-promote block, where our task creates parallelism from the loop induction variable `a`. The parallelism is created by dividing up the remaining iterations into two parts, with half going to the child task, and half to its parent. From this point, both parent and child tasks proceed to work on their respective parts of the remaining iterations, but now start from a new block, namely `loop-par`.

**Parallel tasks.** Now that they are running in parallel, our parent and child tasks have to complete their own local computations, and then combine their local parts of the overall result. To perform their local computations, our parallel tasks execute the `loop-par` block, which performs the same steps as the `loop` block (except for line 23). While our parallel tasks execute, additional heartbeat interrupts may trigger additional promotions, thereby recursively manifesting latent parallelism from the tasks. All heartbeat interrupts from hereon jump to the `loop-par-try-promote` handler block, thereby ensuring that all future tasks share the join record in `jr`. After any one of our parallel tasks completes its local work, the task exits (on line 23) by branching to the `exit-par` block. In this block, the task enters the join-resolution protocol, where all parent and child tasks eventually meet with their partners to combine their local results (in the register `r`).

**Join resolution.** While the program is executing, the TPAL runtime keeps a record of the tree induced by the fork instructions issued by tasks. This bookkeeping tree is used by the join instruction to match each task with its parent or child. When a task issues the join instruction, it stashes its register file in the join record and removes the dependency edge on the join point it shares with its partner in the tree. The first task to complete this step terminates and removes itself from the set of running tasks, and the second performs the next step of join resolution. In this step, the runtime system seeds the task with a new register file. The new register file is obtained by taking the register file of the parent task and extending it with some entries from the register file of the child task. In our `prod` example, this merging process enables the child task to share with its parent the value in its accumulator register `r`. The annotation on the `exit` block

specifies that the contents of register `r` from the child task be copied to register `r2` in the new register file.

After merging register files, the runtime schedules the parent task with this new register file, starting from the *combining block*, which is specified as `comb` in the annotation in our example program. Our combining block takes the sum of the accumulator variables from the parent and child tasks, puts the result in register `r`, and exits by issuing the join instruction again. This time, the join instruction either repeats the process one level up in the tree of tasks, or it reaches the root. If at the root, the join instruction jumps to the original target of the join record, which in our example is the `exit` block. At this point, the parallel program has completed its work.

### 2.3 Nested Parallelism

TPAL can express in a natural way various forms of nested parallelism, for example, in the form of nested loops and recursive functions or their combination. For example, we can implement a “power” function by nesting our running example, `prod`, inside an outer for loop. We parallelize the for loops by using the *outer-loop first* policy, a policy that requires that parallelism is created where it is most beneficial first, that is, from the least recent parallel context. This policy is a necessary condition for any implementation to be backed by the formal efficiency guarantees proved for heartbeat scheduling [5]. As we describe in our evaluation section, TPAL shines in its ability to efficiently execute nested loops, even when their workloads are irregular. The reason is that TPAL can always amortize the cost of parallelism, even across loop boundaries. We present full details of how such a nested parallel programs (including both loops and recursive functions) may be expressed in TPAL in the Appendix.

## 3 Implementation

This section describes the compiler and runtime/OS support needed to translate and execute high-level programs (e.g., Cilk Plus) down to assembly (e.g., x64).

### 3.1 Compiling from C++ to TPAL

To compile a high-level, parallelized loop down to the level of our TPAL assembly, we need to generate the sequential and parallel versions of a loop body, and represent them in the compiler so that the former can be promoted to the latter on-demand. We present our technique by returning to our running `prod` example, this time starting from a Cilk Plus version.

```
void prod_cilk(int a, int b, int* c) {
    reducer_opadd<int> r(0);
    cilk_for (; a != 0; a--) { r += b; }
    *c = r; }
```

The program uses the syntactic parallel-loop extension provided by Cilk Plus to compute the result. Its loop body uses Cilk’s idiomatic pattern for accumulating results, namely

reducer variables [29]. There are similar mechanisms in OpenMP [53] and TBB [41]. The reducer variable enables tasks to work independently on their own local view of  $r$ . When tasks join their results, they sum their local views to compute the final result.

In any such linguistic parallelism, there is necessarily some representation of high-level parallel constructs, e.g., `cilk_for`, `reducer`, that stays intact for some subset of the early stages of the compilation pipeline. Then, at some later stage, the high-level parallel constructs are lowered by that stage into simpler forms. However, the code resulting from lowering pass can easily block optimizations in subsequent passes, such as loop-invariant code motion, loop vectorization, etc. As such, there is motivation to keep the high-level structure and make optimization passes aware of it. This approach is exemplified by recent work on Tapir, an effort to extend LLVM's IR with support for Cilk-style parallelism [56]. Tapir carefully extends LLVM to allow high-level parallel constructs to propagate late into the compilation pipeline, thereby unlocking compiler optimizations that were otherwise inaccessible. Although our TPAL does not yet have compiler support, we believe that a new implementation combining Tapir and TPAL is feasible, and can benefit from the best of both worlds: (1) traditional compiler optimization of high-level parallel code for early stages, thanks to Tapir, and (2) efficient, granularity control, a la TPAL, for later stages. The idea is to implement a new pass that lowers high-level parallel constructs into TPAL instructions. We summarize the lowering steps below.

**Code versioning.** We show in Figure 4 the fragments needed for our `prod` example. The first piece of code is the `prod` function, which represents the initial serial-by-default part of the loop. It corresponds to the first three assembly blocks shown in Figure 2.

The parallel fragment needs to make use of the TPAL runtime system. We show its interface in Figure 3. The interface exports definitions for join-record objects, and `fork` and `join` functions, corresponding to the linguistic forms our TPAL formalism.

The next code fragment, namely `prod_par`, corresponds to the parallel blocks of assembly. It represents one chunk of work, executing sequentially on one processor, but in the context of a parallel execution. After finishing its local work, the function enters the join protocol by calling the `join` function. The final fragments of `prod` are its heartbeat handlers. For convenience, we changed the conventions relating to the `fork` instruction slightly compared to its counterpart in the formalism. Our `fork` instruction takes, in addition to the join record, closures corresponding to the child and parent tasks, and the combine block. Implicit in this convention is that the parent task is rematerialized by the handler function (and correspondingly, the interrupted task exits early to avoid duplicating work).

```
class joinrec { ... };
template <class Child, class Parent, class Comb>
void fork(joinrec* jr, Child c, Parent p, Comb m);
void join(joinrec* jr);
```

**Figure 3.** Scheduling interface used by application code scheduled by TPAL runtime.

```
1 void prod(int a, int b, int* c) {
2   int r = 0;
3   for (; a != 0; a--) { r += b; }
4   *c = r; }
5 void prod_par(int a, int b, int* c, joinrec* jr) {
6   int r = *c;
7   for (; a != 0; a--) { r += b; }
8   *c = r;
9   join(jr); }
10 bool loop_try_promote(int a, int b, int* c) {
11   if (a < 2)
12     return false;
13   joinrec* jr = new joinrec;
14   loop_promote(a, b, c, jr);
15   return true; }
16 void loop_promote(int a, int b, int* c,
17                  joinrec* jr) {
18   int m = a / 2; int n = a % 2;
19   int* rs = new int[2];
20   fork(jr,
21     [=] { // child
22       rs[1] = 0; prod_par(m, b, &rs[1], jr);
23     }, [=] { // parent
24       rs[0] = *c; prod_par(m + n, b, &rs[0], jr);
25     }, [=] { // combine
26       *c = rs[0] + rs[1];
27       delete [] rs; join(jr)
28     }); }
```

**Figure 4.** Code fragments of our `prod` program.

### 3.2 Safe & Efficient Heartbeat Triggering

Next, we describe the compiler, runtime, and OS support to enable the serial version of a loop body to be promoted to our parallel version for the next chunk of iterations.

**Rollforward compilation** In our formal model of TPAL, we assume that heartbeat handlers are triggered every time a task meets two conditions: (1) at least  $\heartsuit$  cycles passed since the previous handler invocation and (2) the control flow enters a promotion-ready program point. Although we cannot rely on any ready-made mechanism from the OS, we can build one on top of OS signaling. There is a challenge, however: given that an interrupt triggered by an OS signal may arrive at *any* step of execution of a task, we somehow need to satisfy the second condition above. In other words, we need a mechanism that triggers a heartbeat interrupt downstream from wherever the interrupted task might currently be executing, given just the current program context. Fortunately, there is a ready-made solution for this problem: we can use

the classic technique of *rollforward compilation* [49]. Rollforward compilation is a general technique for protecting sections of code from being interrupted by OS signals. The idea is to compile a sequence of instructions so that, when preempted by a signal, the sequence executes to its end, and then invokes the signal handler. Our insight is that we can employ rollforward for parallel loops by treating as critical sections the control paths that exist between promotion-ready program points.

We implemented a small rollforward compiler for x64 assembly. The output of this compiler consists of two versions of the input program: original and the rollforward versions. The original version is a copy of the input assembly, modified only so that each line is labeled, e.g., `o0`, `o1`, and so on. As such, there is negligible runtime overhead cost the original program. The rollforward version differs from the original in two respects. First, it has different line labels, e.g., `r0`, `r1`, and so on. Second, any instruction in the rollforward version that jumps to a promotion-ready program point jumps instead to the corresponding handler function. The overall effect is that the original and rollforward instructions align perfectly up to instruction labels, the original version behaves such that it never triggers a heartbeat interrupt, and the rollforward version such that it always triggers a heartbeat interrupt at the next promotion-ready program point.

**Enabling promotions** To enable rollforward at runtime, we need to assign the TPAL worker threads a (Linux) signal handler. When it initializes, the TPAL runtime configures an alarm to invoke this interrupt handler every  $\heartsuit$  microseconds. When it is invoked, the handler uses a table that maps from labels in the source program to labels in its rollforward version. e.g., for x64 prod blocks,  $\{o0 \rightarrow r0, \dots, o7 \rightarrow r7\}$ . This table is generated by the compiler, and is loaded once, by the binary load routine, before the TPAL runtime initializes. When it receives a signal, the handler inspects the program counter of the OS thread that was interrupted by the signal. If the program counter matches a key in the table, the handler replaces that program counter by the corresponding rollforward entry. As a consequence, the program will continue executing until it enters the next promotion-ready program point, at which time it will invoke a handler function.

In general, to make it safe to invoke a handler function, we need support from the compiler to generate *compensation code*. Compensation code is needed because compiler optimization passes may create drift between the program variables and the arguments expected by the handler function. The compensation code materializes the live variables from the registers and stack at the promotion-ready program point, and calls the handler function. For example, suppose we compile a program that iterates over an array, using a pair of integer indices to point to the next cell of the array and the length of the array, respectively. A compiler optimization

may change the types of integer indices used in an array traversal to be direct pointers on the array, thereby creating an incompatibility with a handler function, which may expect as arguments the integer representation of the indices. This problem exists in many other contexts, such as debugging and JIT compilation, and fortunately, there is a general approach for solving it, namely *on-stack replacement (OSR)* [23]. Prior work shows that OSR does not significantly degrade code quality when there are a small number of points in the code that require replacing the stack [27, 39]. This is our case where we only have one promotion-ready program point per parallelized loop.

### 3.3 Taming Code Bloat

Our compilation technique causes an increase in the size of the program binary. Overall, the increase is in linear proportion to the size of the *parallel* regions of code in the program, i.e., blocks of code containing *spawn/sync* calls and *parallel-for* loops. This increase in code size can increase instruction-cache misses, but the misses remain under control, because transfers to the rollforward code are amortized by the heartbeat.

The potential blowup in code size due to the introduction of serial, parallel, and handler blocks can be controlled by compiler support. The handler blocks are likely to impose only marginal cost, whereas the issue of emitting different serial and parallel blocks is a bit more subtle, but nevertheless introduces only a modest tradeoff between the advantage of having different blocks for serial and parallel loop bodies, e.g., *loop* and *loop-par*, versus using one block that is the merging of *loop* and *loop-par*.

### 3.4 Benchmark Implementation

For our benchmark implementations, we used a manual version of on-stack replacement. We use existing compiler mechanisms to generate the compensation code required for each of our benchmark programs. In Figure 5, we demonstrate this technique with our prod program. The first step is to declare a flag, namely *heartbeat*, to stand in for an interrupt delivery. If the conditional branch at line 5 sees true, then the program invokes the handler function. However, we intervene so that this conditional branch is eliminated and therefore never executes at runtime (so we never need to allocate memory for the *heartbeat* global). We simply compile this C++ code to assembly, pass the assembly through our rollforward compiler, and complete the process by manually editing the assembly code generated by rollforward. The edits consist of eliminating from the final assembly all of the conditional branches for the *heartbeat* global. We replace each such conditional in the non-rollforward blocks with *nop* instructions, and we replace each such conditional in the rollforward blocks with an unconditional jump to the handler function. These changes together achieve the desired behavior: the non-rollforward part of the program skips over

```

1 extern volatile bool heartbeat;
2 void prod(int a, int b, int* c) {
3     int r = 0;
4     for (; a != 0; a--) { r += b;
5         if (heartbeat &&
6             (*c = r; loop_try_promote(a, b, c)))
7             return; }
8     *c = r; }

```

**Figure 5.** The instrumented version of our prod program, using our semi-manual rollforward compilation technique.

the conditional, whereas the rollforward version jumps to a certain, compensation block. The compensation block materializes all the program state from the running program, e.g., prod, and calls the corresponding handler function, e.g., loop\_try\_promote. This instrumentation has close to zero cost in the common case, excluding any indirect costs associated with the nop instructions, which can be avoided with compiler support.

**Signaling in Linux.** In Linux, there are at least two off-the-shelf mechanisms we can leverage for heartbeat signals. The first one we call the *ping thread* because it employs a dedicated OS thread to send heartbeat signals to the worker threads. While the approach is simple, the linear signaling does not scale with large core counts, and unfortunately the pthreads library does not offer a broadcast/multicast interface. Some modern architectures expose programmable interrupts based on hardware performance counters. For example, the PAPI library [50] allows for interrupts to be raised on a per-core basis when the cycle counter on the appropriate core exceeds some programmer-specified threshold. Our implementation of TPAL can be configured to use the simple ping-thread approach or the Linux-based PAPI approach. Neither of these interfaces is a natural fit for our needs and neither is particularly optimized for low-latency delivery. Both software and hardware overheads in signal delivery have tangible effects on the achievable heartbeat frequency. In Section 4, we mitigate these effects using custom OS support.

## 4 Evaluation

We compared our TPAL-based Heartbeat Scheduling implementation with the state-of-the-art Cilk Plus system using a common set of benchmarks on a 16 core machine. This comparison is complex because, while Cilk’s execution model performs an initial decomposition when latent parallelism is encountered, TPAL’s execution model involves recurrent decomposition on each beat. It is also important to understand that, beyond the additional overheads incurred by TPAL, the amount of latent parallelism and whether it makes sense to manifest it, varies from benchmark to benchmark. Our goal is to effectively leverage the latent parallelism when it exists and is useful, while paying no cost when it does not exist or is not useful.

Our results show that TPAL creates tasks with considerably lower overhead than Cilk (geomean 13.8× lower), and hence can manage the finer granularity tasks that necessarily result from recurrent decomposition. At the maximum available scale, TPAL achieves significant speedups over Cilk (geomean 53%) for benchmarks that are amenable to recurrent decomposition, while the others (that do not lend well to our technique) incur minimal slowdown (geomean 9.8%).

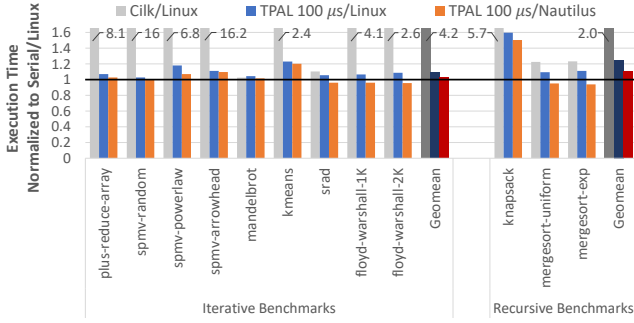
Of course, if TPAL’s overheads were zero, no slowdown would ever occur, and speedups could be enhanced. We next consider TPAL’s overheads in detail. TPAL’s compile-time transformations do not create significant overhead in the transformed code compared to original sequential code. The primary sources of overhead are due to the promotion process, and the heartbeat interrupt mechanism that triggers it. This section uses Linux signals as the mechanism. In the next section, we consider other mechanisms in pursuit of lower-overhead and finer-granularity heartbeat interrupts.

### 4.1 Benchmarks

**Iterative (loop-based) benchmarks.** We ported the *kmeans* and *srad* benchmarks from the Rodinia benchmark suite [21]. For *kmeans*, we use an input of 1 million objects, and for *srad* an 4k × 4k input matrix. The *spmv* benchmark is the classic sparse-matrix by dense-vector product algorithm. Its sparse-matrix input is represented in the compressed sparse row (CSR) format, with non-zero elements represented by double-precision floats. The random matrix is a sparse matrix with 273 million non-zero elements (and non-empty), and a maximum column size of 100. The powerlaw matrix is a random matrix with 186 million non-zero elements, with a power law characteristic [52]. Its largest column contains 5 million non-zero elements, which is 3% of the total number of non-zero elements in the matrix. The arrowhead matrix is a structure that is noted for being particularly challenging for task scheduling [59]: its diagonal, first column, and first row are filled with non-zero elements. The *floyd-warshall* benchmark is a purely loop-based implementation of the classic algorithm for finding the shortest path in a weighted graph. The *mandelbrot* benchmark computes a square image representation of the mandelbrot set [22].

**Recursive benchmarks.** We ported the *knapsack* and *mergesort* benchmarks from the Cilk benchmark suite [44]. The *knapsack* benchmark is the only one of our benchmarks that is non-deterministic: the amount of work it performs depends on the schedule. The *mergesort* algorithm is the only benchmark that uses both parallel loops and parallel recursive calls. In particular, the outer sort function and the inner merge function expose parallelism in a recursive, divide-and-conquer fashion, but there is also a parallel copy operation that moves items to and from a temporary buffer, which exposes parallelism via a parallel loop. The inputs of *mergesort* are generated from uniform and exponential distributions.





**Figure 6.** Task creation overheads for Cilk Plus and TPAL. In contrast to Cilk Plus, TPAL’s overheads are minimal.

## 4.2 Experimental Setup

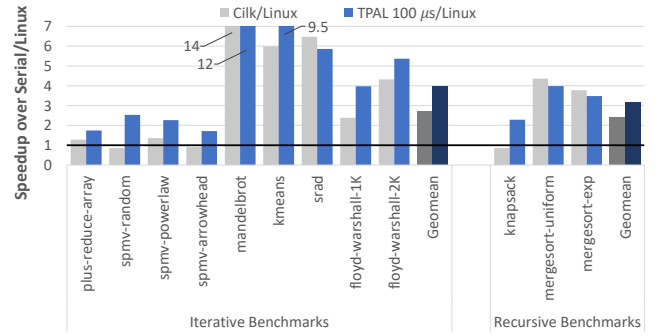
Our primary test bench is a Dell PowerEdge R6415, with a single-socket, 16-core AMD EPYC 7281 (Naples) processor running at 2.7GHz with 64KBL1i, 32KBL1d and 512KB L2 per core, and 32MB of shared L3 cache. It has one NUMA node and 32GB of DDR4 2400MHz RAM. Our machine runs Fedora Server 32, with stock Linux Kernel 5.8.13-200. We disabled SMT (hyperthreading) and we configure the machine’s BIOS to use the maximum performance profile (thus disabling DVFS). For our test machine, we assigned the heartbeat rate to be  $\heartsuit = 100\mu s$  by following the tuning process proposed in the original heartbeat paper [5]. Our code is compiled with GCC version 7.5.0 with flags `-O3 -m64 -march=x86-64`. For load balancing, our TPAL runtime and that of Cilk Plus use randomized work stealing.

We reserve the first core either to do nothing during the benchmark or to execute the ping thread, when needed. The reasons we picked this configuration are (1) we wanted to avoid the overheads generated by a ping thread from affecting any of the worker threads that executed the benchmark workload and (2) we can avoid various other sources of overhead that can, in Linux and Nautilus kernels, slow down the first core in the system.

The Serial/Linux programs are, in all cases but one, the versions of the parallel programs with spawns/joins (also parallel loops) removed. Only in the case of mergesort did we use a significantly different serial program, i.e., serial mergesort. We report the average over 30 runs.

## 4.3 TPAL vs Cilk Plus

**TPAL’s task creation overheads are lower than Cilk’s.** Cilk Plus benefits from significant engineering to make its parallel function call mechanism efficient, even when the program runs on a single core [30]. Furthermore, its parallel loop construct implements an additional form of granularity control by splitting its parallel loop range into  $8P$  blocks, where  $P$  denotes the number of cores. The single-core execution of Cilk Plus is nevertheless slowed down noticeably by task-creation costs in certain cases, whereas in our approach, there is one uniform mechanism, namely the heartbeat, that ensures task-related overheads are well amortized for *all*



**Figure 7.** Speedups over serial execution on Linux, 15 cores. Overall TPAL outperforms Cilk Plus.

*programs and all inputs.* Figure 6 compares the single-core running times of Cilk Plus and TPAL versions of the benchmarks. In all cases but one, our implementations are as fast or faster than those of Cilk Plus. The only exception is *mandelbrot*, which is 2% slower.

**TPAL performs better than Cilk at full scale.** To scale up, both implementations need to create or promote a number of tasks that is sufficient to keep the cores fed, while also keeping task overheads low. Figure 7 compares the 15-core running times of Cilk Plus and TPAL versions of the benchmarks. The parallel execution times of the TPAL versions are considerably lower than those of Cilk, with four exceptions. Of these exceptions, only *mandelbrot* rises above 10%. When TPAL’s recurrent decomposition comes into play, speedups of 53% (geomean) compared to Cilk result, while when recurrent decomposition does not come into play, slowdowns of only 9.8% (geomean) occur.

The *floyd\_warshall* benchmark makes for an interesting point of comparison, because it shows a situation where Cilk’s granularity-control heuristic for `cilk_for` loops fails. For the input size of 1k vertices, there is not enough parallelism to keep all 15 cores fed. The Cilk heuristic generates for this input a much larger number of tasks than our technique (23× more). As a consequence Cilk Plus achieves a higher utilization than TPAL (82% vs. 54%, respectively), and yet is 67% slower (Figure 7), owing to high task overheads. In effect, the Cilk version keeps processors fed doing the busy work of creating and destroying an overabundance of tasks, and ultimately performs worse for it. Our approach finds the right balance for this input, creating just enough tasks to keep cores fed with useful work, but not wasting time sharing too-small tasks. Moreover, as the input increases to 2K vertices, thereby favoring Cilk’s granularity heuristic, our approach scales as well as Cilk does, taking advantage of an amount of parallelism that is closer to keeping all cores fed and reaches comparable utilization levels (89% for Cilk Plus vs. 84% for TPAL). We redirect the interested reader to the Appendix for more details.

Overall, these results show that our approach achieves excellent performance both on parallel as well as sequential runs, i.e., scaling up and down.

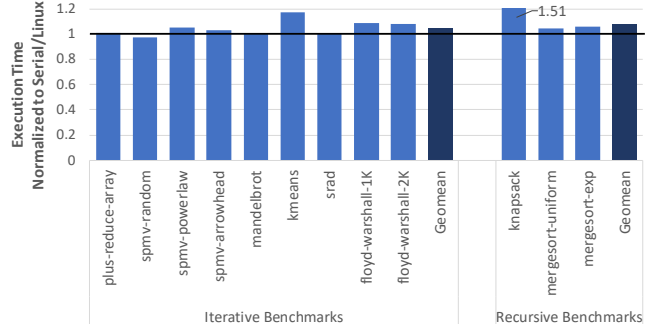
#### 4.4 TPAL Approaches Near-Zero Cost

We now consider the sources of overhead that could limit the performance of TPAL. If these overheads were zero, we would expect the recurrent decomposition of TPAL to always perform strictly better than the initial decomposition of Cilk.

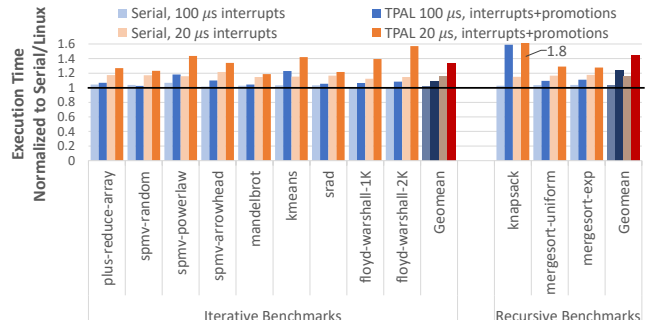
**TPAL’s compilation-related performance overhead is low.** Thanks to its serial-as-the-default scheduling policy, our compilation technique produces binaries whose performance is close to that of corresponding serial programs, because they are very close to the serial programs. Figure 8 compares the execution time of the serial baseline programs of our benchmarks with those of our TPAL binaries. Here, the heartbeat interrupt mechanism is turned off, so pure sequential execution occurs. Our programs are at most 6% slower than their serial baseline programs, except for *floyd-warshall*, *kmeans*, and *knapsack*. The performance of *floyd-warshall* may be related to our manual compilation technique having to slightly modify the innermost loop which is very fine grained, causing a perturbation. In an integrated compiler-based implementation of TPAL (Section 7), that would be easily avoided. The *kmeans* benchmark is slower by 17% because the TPAL version uses an auxiliary data structure to accumulate centroid values, whereas the serial program does not. This situation is the same in the original Rodinia implementations, and as such is not a limitation of our compilation technique.

The 51% slowdown of *knapsack* is the most concerning. It happens for a simple reason: although this benchmark performs almost no computation besides recursive calls, our implementation still pays a cost for pushing and popping promotion-ready stack marks. This cost is visible in *knapsack*, because there is little other computation. In contrast, although it also incurs the costs of maintaining the promotion-ready stack marks, *mergesort* shows only 4–6% overhead, suggesting the bookkeeping costs are less significant. Additional optimizations (e.g., 32-bit pointers) may reduce this overhead but are outside the scope of this work.

**Signal overhead is low, but could be improved.** To promote latent parallelism in Heartbeat Scheduling, signals must interrupt each core on a regular basis. These interrupts need to arrive often enough to create sufficient parallelism, but far enough apart to keep overheads low. Figure 9 isolates and presents the overheads on a single core due to the interrupt mechanism only (i.e., without any promotions), as well as when interrupts generate parallel tasks. The bars labeled Serial represent runs of the *serial baseline* program (not TPAL), and therefore help to isolate the cost of interrupts. The figure presents results only for the INT-PingThread approach of producing the beat, as described in Section 3.2. We



**Figure 8.** Normalized execution time of TPAL sans heartbeat interrupts and concomitant promotions, on Linux, single core. The compilation-related performance overhead of TPAL is minimal compared to sequential code.

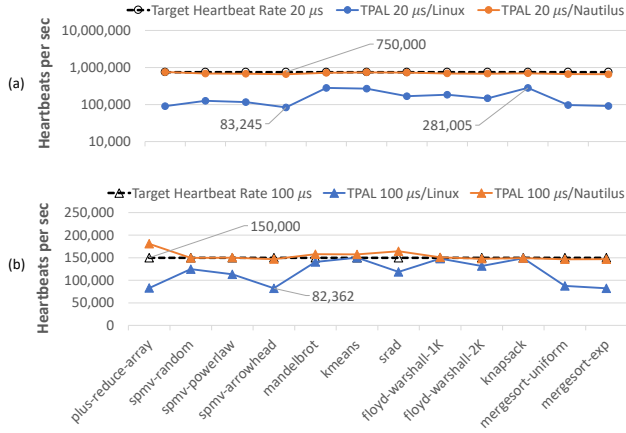


**Figure 9.** TPAL overheads including interrupts only, and interrupts plus promotions, on Linux, single core.

do not present the INT-Papi approach as it always incurs much higher overheads and does not provide any additional benefits.

Interrupt-only overheads at  $\heartsuit = 100 \mu s$  are generally low, with a geomean of 3%. At  $\heartsuit = 20 \mu s$ , however, interrupt-only overheads approach 20% on several occasions, leading to a geomean as high as 16%. As we describe in more detail in Section 5.1, these overheads have a considerably higher and compounding effect at larger scales. However, it is possible to mitigate them through the direct use of timer and IPI hardware as enabled in Nautilus.

**Promotion overhead is low but could be improved.** Figure 9 also depicts the total overhead of promotions by taking the running times of our TPAL programs on a single core when allowing not only signals, but also promotions, to happen. At  $\heartsuit = 100 \mu s$ , all benchmarks except one input of *spmv*, *kmeans*, and *knapsack*, incur a low overhead at or below 11%. The overheads of *kmeans* and *knapsack* are already explained by the compilation-related overheads, which are not specific to TPAL. Although we cannot fully explain the cause of the 18% overhead of *spmv* with the powerlaw input, we have found that on our other test machines the overhead is closer to 10%. For  $\heartsuit = 20 \mu s$ , however, the overhead of promotions becomes unacceptable and reaches a geomean of 34%.



**Figure 10.** Achieved and target heartbeat rate in Linux and Nautilus, 15 cores.

**Linux misses its target heartbeat rate.** An even bigger problem in Linux is that, owing to high signaling overheads, it largely misses its target heartbeat rate. Figure 10 shows that Linux struggles to keep up with its target rate of 150K heartbeats/s across all 15 cores, even at a leisurely  $\tau = 100\mu s$ . While it gets close to desired rate in 3 out of 12 benchmarks in our suite, in all other cases it misses its target, and most often by a lot. It is important to note that we present results for the best Linux mechanism (INT-PingThread); the INT-Papi mechanism performs even worse. The best Linux mechanism cannot even sustain heartbeat signals at a consistent rate for all benchmarks: for some it is as low as 82K/s, significantly lower than the desired rate of 150K/s.

We do not understand in detail why Linux signaling performs so poorly, but our results are in line with other observations of Linux behavior, as described earlier. Conceivably, a suitable Linux kernel mechanism could be designed to improve the situation, and this is a subject for future work.

The situation is aggravated for  $\tau = 20\mu s$  (Figure 10). Here Linux always misses its mark by a factor of 2.7–9 $\times$ : while the target heartbeat rate is a rapid 750K heartbeats/s across all cores, Linux delivers at most 281K, and as low as 83K. Failure to achieve the target rate corresponds to potentially missed opportunities to extract parallelism from the application. We set out to mitigate this shortcoming by exploring alternate mechanisms to deliver interrupts reliably and at low cost.

**Performance at scale.** In Figure 11, we present speedup curves for all benchmarks. Overall, the curves suggest that TPAL and Cilk Plus achieve scaling as cores are added. However, the curves of TPAL usually show low overhead at small core counts, and highest performance at scale. There is one significant exception: *mandelbrot*. In *mandelbrot*, there is a need to spawn a large number of tasks, in order to keep the 15 cores fed, but the speedup curve of Cilk suggests that TPAL is not generating a sufficient number of tasks. The reason

for insufficient tasks is that the signaling mechanism provided by Linux does not support a high enough throughput to meet the needs. In Nautilus, where the signaling mechanism shows better performance at scale, our *mandelbrot* benchmark scales very well, outperforming the Cilk Plus version.

## 5 OS Support for Heartbeat Scheduling

The signaling mechanisms available in Linux were not designed for the purpose of driving heartbeat interrupts at fine granularity, such as  $\tau = 20\mu s$ – $100\mu s$ . As shown shown by others [33], existing software mechanisms in Linux are unable to achieve predictably low latencies for out-of-band event signaling. While the performance of such signaling mechanisms ultimately depends on hardware capabilities, software overheads, a high degree of abstraction, and mismatched interfaces can introduce barriers to signaling performance. To understand the extent of these issues as they pertain to heartbeat signals, we implemented a prototype in a lightweight OS kernel framework called Nautilus that allows us to precisely control the software overheads of event signaling without significant effort.

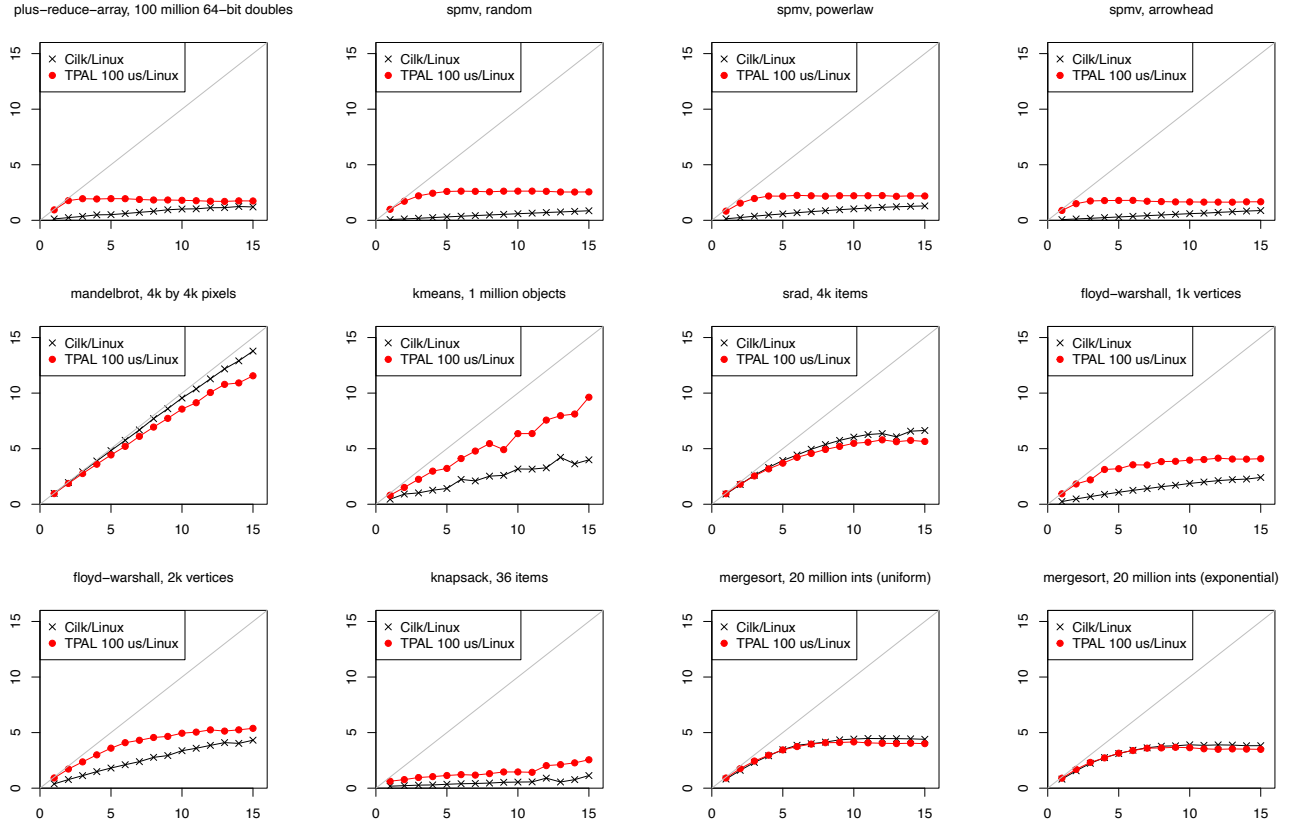
### 5.1 Nautilus and the TPAL HRT

Nautilus is a lightweight kernel framework intended to support *hybrid runtime systems* (HRTs), i.e. language runtimes that have the full power of the OS kernel [34–36]. It is a publicly available open-source codebase<sup>2</sup> that currently runs directly on x64 NUMA hardware and Intel Xeon Phi, as well as in a unikernel configuration atop various virtualization platforms. Applications in Nautilus run in a single shared address space, in kernel-mode, with fully privileged access to the machine.

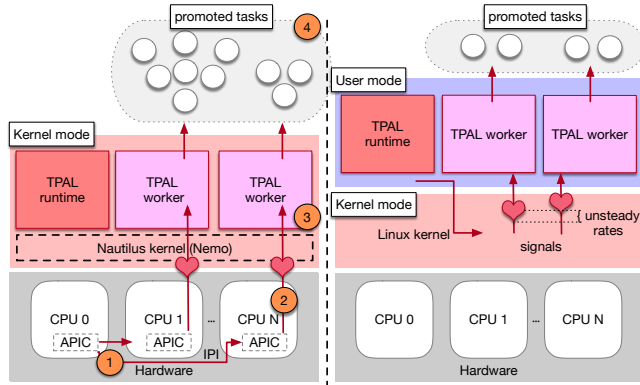
We created a prototype TPAL Hybrid Runtime (TPAL HRT) that runs in Nautilus. It is important to note that TPAL HRT uses application code that is identical to the Linux user-level implementation. The TPAL compiler transformations are no different. The TPAL runtime is also almost entirely identical. What is different is that everything runs within the kernel, and heartbeat signaling is accomplished directly using the x64 timer and interrupt hardware. Arguably, a representation that captures task parallelism, such as TPAL, makes such a radically different implementation feasible for the end-user.

Nautilus includes a lightweight, inter-core signaling framework called Nemo, which reduces signaling latency and jitter significantly [33]. Nemo is essentially a thin veneer around the standard hardware mechanisms for signaling between CPU cores, namely *inter-processor interrupts* (IPIs). In most architectures, IPIs represent the lower limit on architected, out-of-band event signaling. Their cost is typically within a few thousand cycles, most of which is consumed by interrupt handling overhead on the receive side CPU. OSes typically use IPIs for internal synchronization events, but Nautilus

<sup>2</sup><https://github.com/HExSA-Lab/nautilus>



**Figure 11.** Speedup over Serial/Linux for Cilk Plus and TPAL/Linux, varying the number of cores.



**Figure 12.** Heartbeat signaling mechanisms in Nautilus (left) and Linux (right).

exposes this capability to parallel runtimes through Nemo, allowing programmers to multiplex a fixed set of software events and their handler functions atop IPIs. Thus, Nemo is a natural fit for heartbeat signals.

We built our TPAL HRT directly atop Nemo and the hardware timer interrupts that Nautilus also exposes, as depicted in Figure 12. The first TPAL worker on CPU 0 registers a Nautilus timer handler that will be invoked at the specified heartbeat interval ( $\heartsuit$ )/rate. This builds upon the CPU's

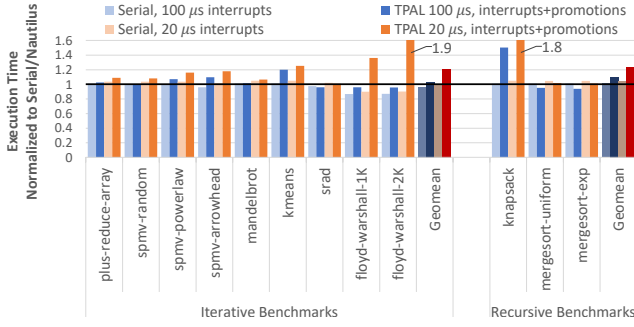
local APIC timer, which can (typically) be programmed to interrupt at intervals down to 10 ns. The timer interrupt directly triggers the heartbeat timer handler. The timer handler in turn uses Nemo to distribute a heartbeat signal to every other core, which Nemo does using IPIs (1). On the destinations, these IPIs trigger (2), via the Nemo framework, the TPAL workers (3), which have more opportunities to unleash parallelism in the form of task promotions (4) due to the consistently achieved heartbeat rate. This approach brings the limit on  $\heartsuit$  closer to the hardware limit.

## 5.2 Signaling Performance in TPAL HRT

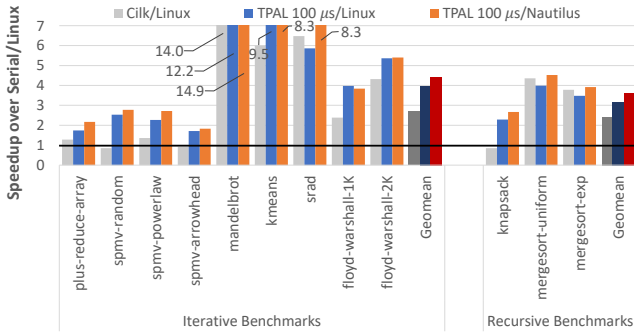
Recall that Figure 9 showed the impact of the heartbeat signaling mechanism on Linux for 20 and 100  $\mu$ s heartbeats. These single-core results were discussed in Section 4.4.

Figure 13 shows the corresponding results for TPAL HRT. The overheads for  $\heartsuit = 100 \mu$ s corresponding to signaling alone are completely masked, whereas in the best-case Linux implementation they were typically around 3–4% and as high as 7%. At 20  $\mu$ s, while the Linux overheads are on the order of 13–22%, the TPAL HRT overheads are at most 4.9%, and are usually much lower. These gains cascade when promotions are also enabled. Clearly, it is Linux that imposes noticeable





**Figure 13.** TPAL overheads including interrupts only, and interrupts plus promotions, on Nautilus, single core.



**Figure 14.** Speedups over serial execution on Linux of 15-core runs of Cilk Plus, TPAL/Linux and TPAL/Nautilus. TPAL outperforms serial/Linux and Cilk Plus on all benchmarks.

costs, even when  $\heartsuit$  is a leisurely 100  $\mu$ s, as current hardware can indeed support lower-cost signaling.

As Nautilus is capable of exploiting the fast interrupt delivery mechanisms of modern hardware, it also achieves its target heartbeat rate. Figures 10 and 13 show that Nautilus practically always achieves the heartbeat rate that it is requested to deliver for both  $\heartsuit = 100 \mu$ s and  $\heartsuit = 20 \mu$ s. While the best Linux mechanism cannot even sustain heartbeat signals at a consistent rate for all benchmarks, even at  $\heartsuit = 100 \mu$ s, Nautilus not only hits the target, but it also delivers a much more consistent rate for both 100  $\mu$ s and 20  $\mu$ s. There is clearly a scaling issue that may affect performance in Linux at higher core counts, but does not affect Nautilus.

### 5.3 Putting It Together: Performance at Scale

Figure 14 compares the overall speedups achieved by Cilk and TPAL on Linux and Nautilus at scale (16 total cores) normalized over the serial execution on Linux. Here  $\heartsuit = 100 \mu$ s, which is generous to Linux.

Cilk Plus achieves speedups on two thirds of our workloads but incurs slowdowns on the remaining third. While some speedups are high (e.g., 14 $\times$  for mandelbrot) often performance improvements are limited, leading to a respectable but less-than-desired speedup geomean of 1.9 $\times$  for iterative benchmarks and 2.4 $\times$  for recursive ones on 15 cores.

TPAL on Linux attains higher performance than Cilk Plus in most cases, or performs comparably well, leading to a

geomean of 4 $\times$  speedup for iterative benchmarks and 3.2 $\times$  for recursive ones. Similarly, TPAL on Nautilus achieves speedup geomeans of 4.4 $\times$  and 3.6 $\times$  respectively for iterative and recursive workloads, outperforming both Cilk Plus and TPAL/Linux in aggregate. Looking at the individual workloads we observe that TPAL on Nautilus achieves the lowest wall-clock execution times than any other system (Cilk plus or TPAL/Linux) for all benchmarks except one: *kmeans*, in which TPAL/Linux narrowly beats TPAL/Nautilus by 12%. It appears that achieving the desired heartbeat interval/rate is a double-edged sword. On the one hand, TPAL HRT can significantly outperform the Linux implementations at scale—for example, in *srab* and *mandelbrot*. Here, the correct and stable heartbeats trigger useful promotions that manifest useful parallelism that contributes to performance, allowing, for example, TPAL to overcome the obstacles it was facing with *mandelbrot* on Linux, as noted in Section 4.3. On the other hand, when Linux fails to achieve the target heartbeat rate, this failure benefits benchmarks in which promotions are *not* desirable by the simple fact that there are fewer promotion opportunities due to this failure.

Overall, TPAL on Linux achieves significant speedups over Cilk (geomean 53%) for benchmarks that are amenable to recurrent decomposition, while the others (that do not lend well to our technique) incur minimal slowdown (geomean 9.8%). It is important to also note that if we consider both Linux and Nautilus implementations, TPAL strictly outperforms Cilk Plus: in all cases, at least one of our TPAL implementations achieves higher speedup than Cilk Plus, and more often both do.

## 6 Related Work

Many task-parallel programming languages have been developed, going back to the 1980s, including multiLisp [37], NESL [12, 14], Cilk (extending C) [30], several extensions of Java [17, 40, 42], parallel Haskell [45, 48, 54], several forms of parallel ML [10, 28, 32, 55, 57, 58, 61], and X10 [20]. All of these task-parallel languages rely on task-scheduling techniques that go back to Brent’s seminal work [18], which has been extended in many directions [4, 8, 11, 13, 15, 16, 24, 43, 51]. These techniques primarily focus on reducing scheduling overheads of tasks that are already created. Task-creation overheads have also proved to be significant, and there has been work on reducing them [6, 25, 30, 38, 41, 60, 63], going back to Cilk-5’s clone optimization. Our TPAL takes inspiration from prior work but takes a different tack: instead of reducing the cost of task creation—which can only be done up to a point—TPAL amortizes the cost against the abundant useful work that a program naturally performs.

There has been recent interest in improving the quality of code generated by high-level parallel languages, such as Cilk Plus. Tapir [56] extends LLVM to support Cilk Plus, bringing to parallel function calls and loops the benefits of LLVM’s existing serial code-optimization passes. Although it

addresses compiler optimization of parallel code, Tapir does not address granularity control, a major challenge in parallelizing codes efficiently, whereas TPAL does. Furthermore, by design, Tapir does not address compiler optimization in the stages following the lowering of Tapir’s high-level parallel instructions, whereas TPAL does. We believe that there is now a feasible path to combine Tapir and TPAL, as we outlined in Section 3.1, which will feature the benefits of each approach in one system.

For our implementation of TPAL, we used the signaling mechanism provided by the OS, in our case, Linux, to drive heartbeats. The performance issues related to the Linux signal mechanism are relatively well known, and are explicitly noted in the source code. In particular, the perf sample rate is limited to 10usec for this reason [2] (lines 418 and 492). There is also some discussion of it in the research literature [62].

Alternatively, the heartbeats can be driven by *software polling* [26]. In software polling, there is typically some sophisticated compiler support that inserts into programs the branch instructions needed to drive heartbeats. This approach has some advantages, especially in the context of managed languages, where there may be preexisting support for software polling. However, there are two challenges facing any implementation of software polling. First, it can block compiler optimizations if not implemented carefully. Second, an implementation needs to ensure there is enough space between polling events to keep polling overheads low, but close enough to consistently meet the target heartbeat rate. These challenges have been addressed by advanced Java runtimes, where there is evidence showing that the overhead cost of the polling is close to 2% [47]. To achieve this result, it is crucial to use a single load/cmp/branch sequence to ensure that the branch would statically be predicted as not taken, and that the register allocation was not affected by the unlikely branch and call [1, 9]. Also, in non-managed languages, such as C++, it can be appropriate to use software polling, and there has been work on bringing compiler support to LLVM [31]. We plan to experiment with such alternative mechanisms in future implementations of TPAL.

## 7 Conclusion

When it comes to performance, it takes two to tango: parallel and sequential computation. Today, we expect the programmer to choreograph carefully, when exactly each will take a step. If parallel goes too far, performance will suffer because of the overheads associated with realizing it in practice, such as task creation, scheduling, etc. If sequential goes too far, scalability will suffer. This choreography involves difficult compromises and requires carefully optimizing code to make sure that each gets their “fair” share by considering everything from lower level concerns such as architectural constant factors, to compiler optimizations, and finally to algorithms.

With widespread availability of multicore hardware, we want languages that encourage and enable writing high-level parallel programs without all of these compromises. To this end, we proposed TPAL as a foundation for writing task-parallel programs. TPAL delivers the right level of parallelism by construction, always and consistently, and applies to irregular, nested, and loop-based parallel codes. We implemented and evaluated TPAL, considering important implementation tradeoffs at the runtime/OS level, on a challenging suite of benchmarks, featuring irregular, fine-grain parallelism. We showed that TPAL achieves consistent efficiency and an ability to find the right amount of parallelism regardless of workload.

## Acknowledgments

**U. Acar and M. Rainey.** This work was partially supported by the National Science Foundation under grant numbers CCF-1901381 and 2028921.

**R. Newton.** This work was made possible by support from the National Science Foundation via awards CCF-2127277 and SPX-1725679.

**P. Dinda and N. Hardavellas.** This work was made possible by support from the National Science Foundation via awards CCF-1533560, CNS-1763743, and CCF-2028851.

**S. Campanoni.** This work was made possible by support from the National Science Foundation via award CCF-1908488 (also, CNS-1763743, and CCF-2028851).

**K. Hale.** This work was made possible by support from the National Science Foundation via awards CNS-1763612, CNS-1718252, and CCF-2028958.

## References

- [1] [n.d.]. Architecture and Policy for Adaptive Optimization in Virtual Machines. <https://researcher.watson.ibm.com/researcher/files/us-groved/RC23429.pdf>. Accessed: 2021-04-1.
- [2] [n.d.]. The Linux source code file core.c. <https://github.com/torvalds/linux/blob/master/kernel/events/core.c>. Accessed: 2021-04-1.
- [3] Umut A. Acar, Vitaly Aksenov, Arthur Charguéraud, and Mike Rainey. 2019. Provably and Practically Efficient Granularity Control. In *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming* (Washington, District of Columbia) (PPoPP '19). 214–228.
- [4] Umut A. Acar, Guy E. Blelloch, and Robert D. Blumofe. 2002. The data locality of work stealing. *Theory of Computing Systems (TOCS)* 35, 3 (2002), 321–347.
- [5] Umut A. Acar, Arthur Charguéraud, Adrien Guatto, Mike Rainey, and Filip Sieczkowski. 2018. Heartbeat Scheduling: Provable Efficiency for Nested Parallelism. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Philadelphia, PA, USA) (PLDI 2018). 769–782.
- [6] Umut A. Acar, Arthur Charguéraud, and Mike Rainey. 2013. Scheduling Parallel Programs by Work Stealing with Private Deques. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (PPoPP '13).
- [7] Umut A. Acar, Arthur Charguéraud, and Mike Rainey. 2016. Oracle-guided scheduling for controlling granularity in implicitly parallel languages. *Journal of Functional Programming (JFP)* 26 (2016), e23.

- [8] Shivali Agarwal, Rajkishore Barik, Dan Bonachea, Vivek Sarkar, R. K. Shyamasundar, and Katherine A. Yelick. 2007. Deadlock-free scheduling of X10 computations with bounded resources. In *SPAA 2007: Proceedings of the 19th Annual ACM Symposium on Parallelism in Algorithms and Architectures, San Diego, California, USA, June 9-11, 2007*. 229–240.
- [9] Matthew Arnold and David Grove. 2005. Collecting and Exploiting High-Accuracy Call Graph Profiles in Virtual Machines. In *Proceedings of the International Symposium on Code Generation and Optimization (CGO '05)*. IEEE Computer Society, USA, 51–62. <https://doi.org/10.1109/CGO.2005.9>
- [10] Jatin Arora, Sam Westrick, and Umut A. Acar. 2021. Provably Space Efficient Parallel Functional Programming. In *Proceedings of the 48th Annual ACM Symposium on Principles of Programming Languages (POPL)*.
- [11] Nimar S. Arora, Robert D. Blumofe, and C. Greg Plaxton. 2001. Thread Scheduling for Multiprogrammed Multiprocessors. *Theory of Computing Systems* 34, 2 (2001), 115–144.
- [12] Guy E. Blelloch. 1996. Programming Parallel Algorithms. *Commun. ACM* 39, 3 (1996), 85–97.
- [13] Guy E. Blelloch, Phillip B. Gibbons, and Yossi Matias. 1999. Provably efficient scheduling for languages with fine-grained parallelism. *J. ACM* 46 (March 1999), 281–321. Issue 2.
- [14] Guy E. Blelloch, Jonathan C. Hardwick, Jay Sipelstein, Marco Zagha, and Siddhartha Chatterjee. 1994. Implementation of a Portable Nested Data-Parallel Language. *J. Parallel Distrib. Comput.* 21, 1 (1994), 4–14.
- [15] Robert D. Blumofe and Charles E. Leiserson. 1998. Space-Efficient Scheduling of Multithreaded Computations. *SIAM J. Comput.* 27, 1 (1998), 202–229.
- [16] Robert D. Blumofe and Charles E. Leiserson. 1999. Scheduling multithreaded computations by work stealing. *J. ACM* 46 (Sept. 1999), 720–748. Issue 5.
- [17] Robert L. Bocchino, Jr., Vikram S. Adve, Danny Dig, Sarita V. Adve, Stephen Heumann, Rakesh Komuravelli, Jeffrey Overbey, Patrick Simons, Hyojin Sung, and Mohsen Vakilian. 2009. A type and effect system for deterministic parallel Java. In *Proceedings of the 24th ACM SIGPLAN conference on Object oriented programming systems languages and applications (Orlando, Florida, USA) (OOPSLA '09)*. 97–116.
- [18] Richard P. Brent. 1974. The parallel evaluation of general arithmetic expressions. *J. ACM* 21, 2 (1974), 201–206.
- [19] F. Warren Burton and M. Ronan Sleep. 1981. Executing functional programs on a virtual tree of processors. In *Functional Programming Languages and Computer Architecture (FPCA '81)*. ACM Press, 187–194.
- [20] Philippe Charles, Christian Grothoff, Vijay Saraswat, Christopher Donawa, Allan Kielstra, Kemal Ebcioglu, Christoph von Praun, and Vivek Sarkar. 2005. 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 (San Diego, CA, USA) (OOPSLA '05)*. ACM, 519–538.
- [21] Shuai Che, Michael Boyer, Jiayuan Meng, David Tarjan, Jeremy W. Sheaffer, Sang-Ha Lee, and Kevin Skadron. 2009. Rodinia: A Benchmark Suite for Heterogeneous Computing. In *Proceedings of the 2009 IEEE International Symposium on Workload Characterization (IISWC '09)*. IEEE Computer Society, USA, 44–54. <https://doi.org/10.1109/IISWC.2009.5306797>
- [22] Intel Corporation. 2014. *Intel C++ Compiler Code Samples*. <https://software.intel.com/en-us/code-samples/intel-compiler/intel-compiler-features/IntelCilkPlus>
- [23] Daniele Cono D'Elia and Camil Demetrescu. 2018. On-stack replacement, distilled. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*. 166–180.
- [24] Derek L. Eager, John Zahorjan, and Edward D. Lazowska. 1989. Speedup versus efficiency in parallel systems. *IEEE Transactions on Computing* 38, 3 (1989), 408–423.
- [25] Karl-Filip Faxén. 2009. Wool-A Work Stealing Library. *SIGARCH Comput. Archit. News* 36, 5 (June 2009), 93–100. <https://doi.org/10.1145/1556444.1556457>
- [26] Marc Feeley. 1993. Polling efficiently on stock hardware. In *Proceedings of the conference on Functional programming languages and computer architecture (Copenhagen, Denmark) (FPCA '93)*. 179–187.
- [27] Stephen J Fink and Feng Qian. 2003. Design, implementation and evaluation of adaptive recompilation with on-stack replacement. In *International Symposium on Code Generation and Optimization, 2003. CGO 2003*. IEEE, 241–252.
- [28] Matthew Fluet, Mike Rainey, John Reppy, and Adam Shaw. 2011. Implicitly threaded parallelism in Manticore. *Journal of Functional Programming* 20, 5-6 (2011), 1–40.
- [29] Matteo Frigo, Pablo Halpern, Charles E. Leiserson, and Stephen Lewin-Berlin. 2009. Reducers and Other Cilk++ Hyperobjects. In *Proceedings of the Twenty-First Annual Symposium on Parallelism in Algorithms and Architectures (Calgary, AB, Canada) (SPAA '09)*. Association for Computing Machinery, New York, NY, USA, 79–90. <https://doi.org/10.1145/1583991.1584017>
- [30] Matteo Frigo, Charles E. Leiserson, and Keith H. Randall. 1998. The Implementation of the Cilk-5 Multithreaded Language. In *PLDI*. 212–223.
- [31] Souradip Ghosh, Michael Cuevas, Simone Campanoni, and Peter Dinda. 2020. Compiler-based timing for extremely fine-grain preemptive parallelism. In *2020 SC20: International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*. IEEE Computer Society, 736–750.
- [32] Adrien Guatto, Sam Westrick, Ram Raghunathan, Umut A. Acar, and Matthew Fluet. 2018. Hierarchical memory management for mutable state. In *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPoPP 2018, Vienna, Austria, February 24-28, 2018*. 81–93.
- [33] Kyle Hale and Peter Dinda. 2018. An Evaluation of Asynchronous Software Events on Modern Hardware. In *2018 IEEE 26th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*. IEEE, 355–368.
- [34] Kyle C. Hale and Peter A. Dinda. 2015. A Case for Transforming Parallel Runtimes Into Operating System Kernels. In *Proceedings of the 24th ACM Symposium on High-performance Parallel and Distributed Computing (HPDC '15)*.
- [35] Kyle C. Hale and Peter A. Dinda. 2016. Enabling Hybrid Parallel Runtimes Through Kernel and Virtualization Support. In *Proceedings of the 12th ACM SIGPLAN/SIGOPS International Conference on Virtual Execution Environments (VEE '16)*. 161–175.
- [36] Kyle C. Hale, Conor Hetland, and Peter A. Dinda. 2016. Automatic Hybridization of Runtime Systems. In *Proceedings of the 25th ACM International Symposium on High-Performance Parallel and Distributed Computing (HPDC '16)*. 137–140.
- [37] Robert H. Halstead, Jr. 1984. Implementation of Multilisp: Lisp on a Multiprocessor. In *Proceedings of the 1984 ACM Symposium on LISP and functional programming (Austin, Texas, United States) (LFP '84)*. ACM, 9–17.
- [38] Tasuku Hiraishi, Masahiro Yasugi, Seiji Umatani, and Taiichi Yuasa. 2009. Backtracking-based load balancing. *Proceedings of the 2009 ACM SIGPLAN Symposium on Principles & Practice of Parallel Programming* 44, 4 (February 2009), 55–64.
- [39] Urs Hölzle, Craig Chambers, and David Ungar. 1992. Debugging optimized code with dynamic deoptimization. In *Proceedings of the ACM SIGPLAN 1992 conference on Programming language design and implementation*. 32–43.
- [40] Shams Mahmood Imam and Vivek Sarkar. 2014. 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*. 75–86.



- [41] Intel. 2011. Intel Threading Building Blocks. <https://www.threadingbuildingblocks.org/>.
- [42] Doug Lea. 2000. A Java fork/join framework. In *Proceedings of the ACM 2000 conference on Java Grande* (San Francisco, California, USA) (JAVA '00). 36–43.
- [43] I-Ting Angelina Lee, Charles E. Leiserson, Tao B. Schardl, Zhunping Zhang, and Jim Sukha. 2015. On-the-Fly Pipeline Parallelism. *TOPC* 2, 3 (2015), 17:1–17:42.
- [44] Charles Leiserson and Aske Plaatt. 1998. Programming parallel applications in Cilk. *SINEWS: SIAM News* 31, 4 (1998), 6–7.
- [45] Peng Li, Simon Marlow, Simon L. Peyton Jones, and Andrew P. Tolmach. 2007. Lightweight concurrency primitives for GHC. In *Proceedings of the ACM SIGPLAN Workshop on Haskell, Haskell 2007, Freiburg, Germany, September 30, 2007*. 107–118.
- [46] Vasileios Liaskovitis, Shimin Chen, Phillip B. Gibbons, Anastassia Ailamaki, Guy E. Blelloch, Babak Falsafi, Limor Fix, Nikos Hardavellas, Michael Kozuch, Todd C. Mowry, and Chris Wilkerson. 2006. Parallel Depth First vs. Work Stealing Schedulers on CMP Architectures. In *Proceedings of the Eighteenth Annual ACM Symposium on Parallelism in Algorithms and Architectures* (Cambridge, Massachusetts, USA) (SPAA '06). Association for Computing Machinery, New York, NY, USA, 330. <https://doi.org/10.1145/1148109.1148167>
- [47] Yi Lin, Kunshan Wang, Stephen M. Blackburn, Antony L. Hosking, and Michael Norrish. 2015. Stop and Go: Understanding Yield-point Behavior. In *Proceedings of the 2015 International Symposium on Memory Management* (Portland, OR, USA) (ISMM '15). Association for Computing Machinery, New York, NY, USA, 70–80. <https://doi.org/10.1145/2754169.2754187>
- [48] Simon Marlow and Simon L. Peyton Jones. 2011. Multicore garbage collection with local heaps. In *Proceedings of the 10th International Symposium on Memory Management, ISMM 2011, San Jose, CA, USA, June 04 - 05, 2011*, Hans-Juergen Boehm and David F. Bacon (Eds.). ACM, 21–32.
- [49] David Mosberger, Peter Druschel, and Larry L Peterson. 1996. Implementing atomic sequences on uniprocessors using rollforward. *Software: Practice and Experience* 26, 1 (1996), 1–23.
- [50] Philip J Mucci, Shirley Browne, Christine Deane, and George Ho. 1999. PAPI: A portable interface to hardware performance counters. In *Proceedings of the department of defense HPCMP users group conference*, Vol. 710.
- [51] Girija J. Narlikar and Guy E. Blelloch. 1999. Space-Efficient Scheduling of Nested Parallelism. *ACM Transactions on Programming Languages and Systems* 21 (1999).
- [52] MEJ Newman. 2005. Power laws, Pareto distributions and Zipf's law. *Contemporary Physics* 46, 5 (Sep 2005), 323–351. <https://doi.org/10.1080/00107510500052444>
- [53] OpenMP Architecture Review Board. [n.d.]. OpenMP Application Program Interface. <http://www.openmp.org/>
- [54] Simon L. Peyton Jones, Roman Leshchinskiy, Gabriele Keller, and Manuel M. T. Chakravarty. 2008. Harnessing the Multicores: Nested Data Parallelism in Haskell. In *FSTTCS*. 383–414.
- [55] Ram Raghunathan, Stefan K. Muller, Umut A. Acar, and Guy Blelloch. 2016. Hierarchical Memory Management for Parallel Programs. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming* (Nara, Japan) (ICFP 2016). ACM, New York, NY, USA, 392–406.
- [56] Tao B. Schardl, William S. Moses, and Charles E. Leiserson. 2017. Tapir: Embedding Fork-Join Parallelism into LLVM's Intermediate Representation. In *Proceedings of the 22nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (Austin, Texas, USA) (PPoPP '17). Association for Computing Machinery, New York, NY, USA, 249–265. <https://doi.org/10.1145/3018743.3018758>
- [57] K. C. Sivaramakrishnan, Lukasz Ziarek, and Suresh Jagannathan. 2014. MultiMLton: A multicore-aware runtime for standard ML. *Journal of Functional Programming* FirstView (6 2014), 1–62.
- [58] Daniel Spoonhower. 2009. *Scheduling Deterministic Parallel Programs*. Ph.D. Dissertation. Carnegie Mellon University. <https://www.cs.cmu.edu/~rwh/theses/spoonhower.pdf>
- [59] Torbjørn Johnsen Tessem. 2013. *Improving parallel sparse matrix-vector multiplication*. Master's thesis. The University of Bergen.
- [60] Alexandros Tzannes, George C. Caragea, Uzi Vishkin, and Rajeev Barua. 2014. Lazy Scheduling: A Runtime Adaptive Scheduler for Declarative Parallelism. *TOPLAS* 36, 3, Article 10 (Sept. 2014), 51 pages.
- [61] Sam Westrick, Rohan Yadav, Matthew Fluet, and Umut A. Acar. 2020. Disentanglement in Nested-Parallel Programs. In *Proceedings of the 47th Annual ACM Symposium on Principles of Programming Languages (POPL)*.
- [62] Xi Yang, Stephen M. Blackburn, and Kathryn S. McKinley. 2015. Computer Performance Microscopy with Shim. *SIGARCH Comput. Archit. News* 43, 3S (June 2015), 170–184. <https://doi.org/10.1145/2872887.2750401>
- [63] Christopher S Zakian, Timothy AK Zakian, Abhishek Kulkarni, Budhika Chamith, and Ryan R Newton. 2015. Concurrent Cilk: Lazy Promotion from Tasks to Threads in C/C++. In *International Workshop on Languages and Compilers for Parallel Computing*. Springer International Publishing, 73–90.