

Locality-Aware Dynamic Task Graph Scheduling

Jordyn Maglalang*, Sriram Krishnamoorthy[†], Kunal Agrawal*

*Washington University in St. Louis

[†]Pacific Northwest National Laboratory

jordyn.maglalang@wustl.edu, sriram@pnl.gov, kunal@seas.wustl.edu

Abstract—Dynamic task graph schedulers automatically balance work across processor cores by scheduling tasks among available threads while preserving dependences. In this paper, we design NABBITC, a provably efficient dynamic task graph scheduler that accounts for data locality on NUMA systems. NABBITC allows users to assign a color to each task representing the location (e.g., a processor core) that has the most efficient access to data needed during that node’s execution. NABBITC then automatically adjusts the scheduling so as to preferentially execute each node at the location that matches its color—leading to better locality because the node is likely to make local rather than remote accesses. At the same time, NABBITC tries to optimize load balance and not add too much overhead compared to the vanilla NABBIT scheduler that does not consider locality. We provide a theoretical analysis that shows that NABBITC does not asymptotically impact the scalability of NABBIT.

We evaluated the performance of NABBITC on a suite of benchmarks, including both memory and compute intensive applications. Our experiments indicate that adding locality awareness has a considerable performance advantage compared to the vanilla NABBIT scheduler. Furthermore, we compared NABBITC to both OpenMP tasks and OpenMP loops. For regular applications, OpenMP loops can achieve perfect locality and perfect load balance statically. For these benchmarks, NABBITC has a small performance penalty compared to OpenMP due to its dynamic scheduling strategy. Similarly, for compute intensive applications with course-grained tasks, OpenMP task’s centralized scheduler provides the best performance. However, we find that NABBITC provides a good trade-off between data locality and load balance; on memory intensive jobs, it consistently outperforms OpenMP tasks while for irregular jobs where load balancing is important, it outperforms OpenMP loops. Therefore, NABBITC combines the benefits of locality-aware scheduling for regular, memory intensive, applications (the forte of static schedulers such as those in OpenMP) and dynamically adapting to load imbalance in irregular applications (the forte of dynamic schedulers such as Cilk Plus, TBB, and Nabbit).

I. INTRODUCTION

In recent years, parallel computers have become ubiquitous and many high-level programming languages and libraries, such as OpenMP, Cilk Plus, Nabbit, etc. have emerged. These languages and libraries allow programmers to express the logical parallelism in their programs while the runtime scheduler schedules the work on the available cores automatically. For multicores with few cores and uniform access to the memory hierarchy, these languages and runtime systems provide both good performance and a relatively simple programming model.

On large multicores with non-uniform memory access (NUMA), however, *locality* is an important consideration — a *remote memory access*—access to data reachable from a memory controller that is further away via the on-chip network—can cost much more than a *local memory access*. Regular applications can be structured to implicitly ensure

locality between initialization and subsequent use when using static schedulers such as in OpenMP. Irregular applications, on the other hand, require dynamic load balancing which dynamic schedulers, such as those in OpenMP tasks, TBB, and Cilk Plus provide. These systems generally have no notion of the location of data and often fail to provide good performance for regular memory-intensive applications.

Ideally, one would like to have a high-level and easy-to-use programming model which incorporates dynamic scheduling and locality. We present NABBITC, a locality-aware extension of a task graph library NABBIT. In the NABBIT programming model, the programmer expresses computations as a task graph where each node is a task and edges represent dependences between tasks. NABBIT is a library built on top of a Cilk Plus¹ and therefore, NABBIT programs are scheduled using a provably good work-stealing scheduler. This paper makes NABBIT locality aware by allowing the programmer to give locality hints to the scheduler using a simple coloring scheme. In particular, we make the following contributions.

1. We extend the interface so that the programmer can provide a *color* to each task; if a task is colored a color c , then the data used by this task is local to processor with color c . Multiple nearby cores can have the same color.
2. We modify both the NABBIT library and the Cilk Plus runtime system to allow processors to preferentially execute tasks that share the color with them. Therefore, if the user provides a “correct coloring”, then workers preferentially execute tasks that access local data, thereby reducing the expensive remote accesses.
3. NABBITC tries to strike a balance between improving locality and preserving the guarantees of low overhead and good load balance provided by NABBIT. We prove that NABBITC, by and large, preserves the asymptotic guarantees provided by NABBIT. In particular, for reasonable task graphs—those with enough parallelism and where tasks of all colors appear near the root of the graph—NABBITC provides speedup that asymptotically is nearly optimal
4. We evaluated the performance of NABBITC on a suite of memory and compute intensive applications and found that it succeeds in providing both good load balance and good locality. Our implementations utilize simple coloring functions, which we detail later, that require minimal additional effort. NABBITC consistently out performs vanilla NABBIT due to improvements in locality. On regular, memory intensive, applications for which OpenMP static schedules are

¹It was originally designed to be built on top of Cilk++, but it is trivial to port to Cilk Plus. Indeed, it was specifically designed such that it can be ported to any programming language that supports fork-join parallelism.

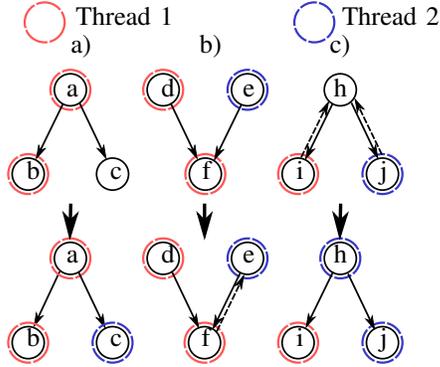


Fig. 1: NABBIT scheduling examples presented in two stages. A thread surrounding a node denotes that thread is processing that node. A solid line from a to b denotes that b is a predecessor of a while a dashed line from b to a denotes that a is in b 's successor list.

best, NABBITC remains very competitive, despite its added overhead of dynamic scheduling. On compute intensive applications, OpenMP tasks edges out as best by exposing more parallelism than OpenMP loops and having a simpler scheduler than NABBIT and NABBITC— however, NABBITC is still competitive in performance. For PageRank, an exemplar irregular benchmark, NABBITC outperforms all other schedulers by combining dynamic load balancing and locality awareness.

II. BACKGROUND

In this section, we describe NABBIT, a high-level task-graph scheduling library built on top of the Cilk Plus runtime system. We outline the NABBIT programming model and show how NABBIT recursively executes task graphs in parallel. We also provide a brief overview of the GCC Cilk Plus implementation upon which NABBIT is built.

NABBIT task-graph scheduling. A NABBIT task graph is a directed acyclic graph with a set of explicit nodes that represent tasks and edges that represent dependences between tasks. Each node u in the task graph specifies its predecessors—tasks that have edges to u and therefore must be executed before u can be processed. For this paper, we will use the terms node and task interchangeably.

We summarize the key aspects of the NABBIT dynamic task graph scheduler (more details in Agrawal et al. [1]). One interesting property of NABBIT is that it computes nodes *on demand*. The scheduler takes an input specified in the form of a sink node, whose execution completes the execution of the task graph. Upon initialization, this node has a list of predecessors and no successors. The sink node together with the predecessor specification transitively identifies all vertices that need to be executed to compute the sink node. The scheduler actions:

- 1) To process a node, a thread initializes the node and its list of predecessors and proceeds to execute them in a recursive parallel depth-first fashion. Consider the example in Figure 1a. When thread 1 wants to process a and finds that b and c are its predecessors that have not been initialized, it goes ahead and tries to process one (b in this case). While b is being processed, another thread can steal c .

- 2) When processing a node's predecessor, if a thread finds that some predecessor has already been initialized by some other thread but has not finished executing, the thread adds the current node to the predecessor's successor list and moves on. In Figure 1b, thread 1 is processing d while thread 2 is processing e . f is a predecessor of both d and e . Each thread will try to initialize the predecessor f but only one will succeed, in this case thread 1. Thread 2 is attaches e to f 's successor list and tries to find other work to do.
- 3) After a node is computed, the thread checks if there are any enqueued successors and if so, determines if those successors are ready to execute (i.e., have no other predecessors on which they are waiting). In the event that a successor is ready, the thread will recursively execute that node. In Figure 1c, both nodes i and j have h enqueued in their successor lists. Thread 1 computes i and checks if h is ready to execute. Since h still depends on j , thread 1 moves on. Thread 2, after computing j checks h , sees that it is ready and proceeds to execute it.

This procedure ensures that a node is computed only after all its (transitive) predecessors have been computed, ensuring correctness. In addition until an initialized node u is computed, it is (a) either in a thread's stack, (b) in a successor list, or (c) is a predecessor of an initialized node. This guarantees that every node u will be inspected and executed eventually. Also, this ensures that the sink node, and thus whole task graph, is executed to completion.

Atomicity choices ensure the absence of data races. The predecessor and successor lists allow threads to execute without blocking/waiting for actions by other threads. The recursive parallel design allows for the implementation of the NABBIT's scheduler as a Cilk Plus program. All vertices in either the predecessor or successor lists can be executed in parallel. In addition, NABBIT ensures that no ordering constraints other than those implied by the predecessor relationships is imposed on the execution: a node u is ready to execute immediately after all its predecessors have been computed and unless every processor is busy doing other work, some processor will find and execute u . This ensures that NABBIT does not alter the task graph's critical path length, enabling the scheduler to guarantee asymptotic optimality. Essentially, if the task graph itself has a parallelism of at least P , then NABBIT guarantees that it gets $\Omega(P)$ speedup on P processors for most reasonable task graphs. In addition, since it leverages the Cilk Plus work-stealing scheduler and uses distributed processing, NABBIT provides low overheads. These properties of asymptotic optimality and low overheads are not normally achieved by other task graph schedulers, such as scheduler currently used in OpenMP's SMPSSs [2], since they do not process nodes on demand.

GCC Cilk Plus We compile NABBIT using the GCC implementation of Cilk Plus, an extension to C++. Cilk Plus is a *processor oblivious* language—the programmer expresses the *logical parallelism* of the program using three keywords without any reference to how many threads must execute the program and how. The `cilk_spawn` keyword indicates that the succeeding function can execute in parallel with its continuation. The `cilk_sync` keyword is a local barrier; all previously

```

1 class DynamicNabbitNode:
2   Key key //this node's key
3   List<Key> predecessors //predecessors' keys
4   List<DynamicNabbitNode *> successors //successor nodes
5   virtual void init() //initialize this code
6   virtual void compute() //computation for this task
7   int color() = color(node.key) //helper function: this node's color

```

9 int color(Key key) // function mapping keys to color

Fig. 2: NABBITC abstract class interface

spawned functions by current function must complete before the program execution can move this statement. Cilk Plus also provides a `parallel_for` keyword, which indicates that all iterations can be executed in parallel. This keyword is essentially syntactic sugar and is implemented using `spawns` and `syncs`.

The Cilk Plus runtime system uses randomized work stealing to schedule these fork-join programs on P available cores. The program executes on P worker threads, one for each core in the target machine. Each worker has a local deque of work. When a worker p executing function `foo` spawns a function `bar`, the frame corresponding to the caller `foo` is placed at the bottom of the p 's deque and p starts executing `bar`. When p returns from a function, it pops the function at the bottom of its deque and continues executing. (If executing on one thread, the program follows the normal depth-first execution followed by C or C++.) If worker q 's deque is empty, it becomes a thief, picks a random victim worker, say p , steals the top frame from p 's deque, and starts executing it. If a steal attempt is unsuccessful, meaning that the victim had an empty deque, then the thief continues to steal until it finds work. The Cilk Plus compiler inserts code at `spawns` and `syncs` to ensure that deques are managed correctly. In addition, when a worker's deque is empty, it makes calls into the runtime to make sure that steals occur correctly.

III. NABBITC DESIGN

In this section, we describe the NABBITC interface and present the extensions to NABBIT and the underlying Cilk Plus runtime that together constitute the NABBITC infrastructure which exploits user provided hints to optimize locality.

NABBITC interface

Figure 2 shows the abstract class interface for defining nodes (tasks) and their data dependencies (edges) in NABBIT. Each node in a task graph inherits from `DynamicNabbitNode` class and is associated with a unique key. Nodes implement the member functions shown in Figure 2. The `init()` and `compute()` functions serve to initialize node parameters and perform the computation represented by the node, respectively. The user specifies the list of the node's predecessors, identified by their keys, in the `predecessors` array. In addition to the information on tasks and their dependencies needed by original NABBIT, NABBITC requires the user to define a `color()` function that returns a node's color. This function definition serves as the mechanism for the user to provide locality information to NABBITC and is the only additional piece of information the user must provide.

```

1 //Recursively spawn colors using morphing continuations
2 void spawn_colors(colors):
3   if len(colors)==1:
4     spawn_nodes(colors[0])
5   else
6     c_p = /*this worker's color*/
7     /*split available colors into two halves*/
8     first_half = colors[0:len(colors)/2]
9     second_half = colors[len(colors)/2:]
10    if c_p in second_half.keys():
11      swap(first_half, second_half)
12    cilkrts_set_next_colors(second_half.keys())
13    cilk_spawn spawn_colors(first_half)
14    spawn_colors(second_half)
15    cilk_sync

```

```

17 //Recursively spawn nodes of the same color
18 void spawn_nodes(nodes):
19   if len(nodes)==1:
20     if nodes[0] is a successor:
21       if nodes[0] is ready:
22         nodes[0].compute_and_notify()
23     else: /*predecessor key*/
24       try_init_compute(this,nodes[0])
25   else:
26     color = nodes[0].color() /*all nodes have same color*/
27     first_half = nodes[0:len(nodes)/2]
28     second_half = nodes[len(nodes)/2:len(nodes)]
29     cilkrts_set_next_colors(color)
30     cilk_spawn spawn_nodes(first_half)
31     spawn_nodes(second_half)
32     cilk_sync

```

Fig. 3: Pseudo-code for color-aware spawning of a set of nodes in NABBITC using morphing continuations. We use a hybrid C++/Python syntax to enhance readability.

Designing a locality-guided task-graph scheduler

NABBITC attempts to achieve multiple goals during scheduling: (1) improve data locality by executing nodes of the same color as the executing processor; (2) achieve good load-balance for the computation as a whole; and (3) introduce minimal overhead into the original NABBIT scheduling pathway. Extending the NABBIT task-graph scheduling library to make use of user-provided locality information involves altering how Cilk Plus workers find and determine what to work on. We introduce two specific changes to NABBIT in order to implement this change in policy: (1) **color-aware scheduling using morphing continuations** allow workers to reorganize work so that they may preferentially execute nodes that have the same color as theirs and (2) **colored steals** allow Cilk Plus workers to find work of their color from the current set of stealable frames. In order to implement these policies, we must also change the Cilk Plus runtime system. We now describe these changes.

Color-aware execution order using morphing continuations. The primary source of concurrency in NABBIT is the concurrent processing of all predecessors or successors² of a given node. As explained in Section II, NABBIT enables this by spawning the execution of all predecessors (or successors) in parallel using a `parallel_for` loop. NABBIT is oblivious to the order in which these nodes are processed. NABBITC, however, extracts colors from this list of nodes in order to preferentially

²As described in Section II, NABBIT implicitly maintains `successors` array for each node u and NABBIT may push successor nodes into it when they must wait for u to complete.

```

1  /*Helper functions to obtain colors*/
2  int color(DynamicNabbitNode node):
3      return node.color()
4  int color(DynamicNabbitNode *node):
5      return node->color()

7  /*Gather list of spawns based on their color.
8   T = Key (for predecessor list) or
9   T = DynamicNabbitNode (for successor list)*/
10 auto gather_colors(T nodes):
11     //group nodes based on their colors
12     Map<int,List<T>> colors
13     for n in nodes:
14         colors[color(n)].add(n)
15     return colors

17 /*Initialize this (already created) node and compute*/
18 void init_node_and_compute():
19     this.init()
20     colors = gather_colors(this.predecessors)
21     spawn_colors(colors)
22     if all this.predecessors have been computed:
23         this.compute_and_notify()

25 /*Try to initialize node's predecessor with key pkey */
26 void try_init_compute(node, pkey):
27     //atomically attempt to create a predecessor with key pkey
28     pred = /*reference to node for key pkey*/
29     if /*creation succeeded*/:
30         pred.init_node_and_compute()
31     else: /*already created by this or some other thread*/
32         atomic pred.successors.add(node) //enqueue

34 /*compute a node and notify its successors*/
35 void compute_and_notify():
36     this.compute()
37     while /*there are new successors in this.successors*/:
38         colors = gather_colors(this.successors)
39         spawn_colors(colors)

```

Fig. 4: Key routines to spawn predecessors and successors in NABBITC.

process nodes with colors that would improve locality.

The crucial function for this purpose is the function `spawn_colors` shown in Figure 3. This function is called on a list of colors colors (and implicitly, a set of nodes which have these colors). At a high-level, when a processor with color c_p is executing this function, it tries to execute the nodes with color c_p by recursively calling `spawn_colors` on the half of the list that contains c_p . Once it reaches the base case (the set colors contains only one color), it then spawns all the nodes of color c_p using the function `spawn_nodes`. This function `spawn_nodes` is essentially a parallel-for loop over the nodes of this color.

The function `spawn_colors` re-organizes the order in which nodes are spawned so that the nodes of the preferred color c_p are spawned first, implementing what we call a *morphing continuation*. The particular strand that is spawned and the continuation of the strand depends on the color of the processor which is doing the spawn. Another important thing to note about this code is that if the preferred color c_p is not present in the list, the function will spawn the nodes in the original ordering of the list — therefore, a worker does not stall even if it can not find the work of its color.

The function `spawn_colors` is called in three places in the NABBITC library. Figure 4 shows the actions to initialize and execute a node. `init_and_compute()` acquires the colors of the current node's predecessors and invokes `spawn_colors()` if there

exist more than one. Similarly, when spawning the list of successors, `compute_and_notify()` collects the set of colors for the list of successors and invokes `spawn_colors()` if there are more than one. Finally, `spawn_colors` is a recursive function which is also called by itself.

This morphing continuation design allows us to use the same mechanism in two scenarios. First, when a processor spawns the predecessors (or successors) if the node it is currently working on, it uses `spawn_nodes` to preferentially execute the predecessor(s) (or successor(s)) of its own color. Second, and the more subtle point, is as follows. Note that in Cilk Plus, when a thief worker steals a task after the spawn of a function³, it executes the function's continuation. Since `spawn_colors` is recursive, when a worker steals a continuation, the first statement it executes is `spawn_colors`. Therefore, the thief also preferentially executes the nodes of *its* color using the same mechanism.

Colored Steals: When a worker has no assigned work, either because it has run out of local work or is at the start of execution, we want that worker to acquire work of its preferred color if possible. In order to do so, we changed the stealing policy of Cilk Plus to allow *colored steal* where a worker checks a deque and only steals the work (continuation) at the top of the deque if that continuation contains some node of this worker's preferred color. We will describe the implementation below — we first describe the policy details about when NABBITC performs colored steals vs. random steals.

One of the goals of NABBITC is to strike a reasonable balance between locality — workers preferentially execute work of their color — and load balance — workers are not idle for too long. In order to do so, we make two changes to the standard Cilk Plus policy of random steals. First, when a worker p with color c_p runs out of work, it does a constant number of colored steal attempts before attempting a random steal. That is, it randomly picks a victim worker q and checks if the frame on the top of q 's deque has any tasks of color c_p — if so, it steals this frame making this a successful colored steal. If not, it tries again. If it fails on a constant number of colored steal attempts, it reverts to a color agnostic random steal. This policy makes sure that p tries to find work of its own color, but then also maintains provable load balance guarantees (as shown in Section IV) by greedily doing any work available if it can not easily find work of its color.

There is an exception to this policy, however. At the beginning of the computation, one worker starts out with executing the root node and all other workers are stealing. At this time, if a worker begins execution in a region of a task graph with no tasks of its preferred color, it will continue executing the available non-preferred tasks until all work is exhausted (as explained in the morphing continuations section). In addition, often, the first steal represents a significant amount of work (conceptually corresponding to nodes higher up in the task graph or computation tree) and a random first steal can potential to lead poor locality. Therefore, we enforce that the first steal a worker performs is a successful colored steal. After the first steal, the worker follows the policy explained above.

³A task is represented at runtime by the task's stack or activation frame

This enforcement does affect Cilk’s time bound, which we explore in Section IV. In our experiments, we found that if all colors are available at the root of the task graph, this time to first work (successful steal) is agnostic to the application, is strictly determined by the number of processors, and, in general, has a small impact on the overall execution time.

We now describe the changes made to both NABBIT and Cilk Plus in order to implement the colored steal policy.

Color-aware GCC Cilk Plus runtime

We make the GCC Cilk Plus runtime color aware by making the following changes. First, we add two additional functions to the Cilk Plus API, shown in Figure 5, that allows NABBITC to provide color information to the runtime system. The first function is straight forward and is simply used by each worker to set the color of this worker. We pin worker threads and assign them a unique color based on their thread id. The second one is used to implement colored steals and requires more explanation. Recall that in order to do colored steals, a thief worker must be able to tell which color nodes are available in the frame that is on the top of victim worker’s deque. This API allows NABBITC to pass this information to Cilk Plus runtime. In particular, before every `cilk_spawn`, NABBITC calls `cilkrts_set_next_colors()` to inform the Cilk Plus runtime about which colors are available in the continuation.

The Cilk Plus runtime is also changed with respect to what it does on spawns. At each spawn, the vanilla Cilk Plus pushes the frame of the currently executing function into a worker’s deque — allowing some other worker to steal the continuation of the spawn. To enable colored steals, we maintain a color deque alongside the work deque to hold the colors available in each continuation. When NABBITC calls `cilkrts_set_next_colors` with a set of colors before the spawn statement, this set of colors is pushed at the bottom of the color deque — therefore, each continuation on the work deque has a corresponding set of colors on the color deque.

Now it is easy to see how one can implement colored steals. When a worker p wants to do a colored steal, it simply checks to see if color c_p (p ’s preferred color) is in the set of colors on the top of victim’s color deque. If so, it pops the top of both the color deque and the work deque and puts them on the top of its corresponding deques making it a successful colored steal. Since the number of colors is determined by the number of workers, we make each entry in the colored deque a fixed length array of boolean flags indicating colors contained in the corresponding continuation. This makes the thief’s check a constant time operation.

Setting continuation colors in NABBITC: As mentioned above, NABBITC must set colors of continuation at each spawn using `cilkrts_set_next_colors` function. This is done on Lines 12, 29, etc within the code in Figure 3. Note that this fits in seamlessly with the design of morphing continuations. At each spawn, we know exactly which colors are available within the spawn and which are available within the continuation. Therefore, NABBITC can easily notify Cilk Plus of the colors available in the continuation by simply telling it which colors are available in the second call to `spawn_colors`.

```
1 void cilkrts_set_worker_color(int color)
2 void cilkrts_set_next_colors(List<int> colors)
```

Fig. 5: Extensions to the Cilk Plus RTS API to inform the runtime of the worker’s preferred color and the colors available in a continuation.

```
1 // Global application parameters
2 int ncolors; // The number of colors being used
3 int n; // The height of the stencil data
4 int B; // The block size of each node
5 int nB; // The number of blocks along the height

7 int color(Key key):
8     int row_id, col_id = get_rowcol_from_key(key)
9     return floor( col_id / nB * ncolors)
```

Fig. 6: An example color function for a stencil application who’s data has been is stored in row-major order. The columns have been distributed evenly, allowing us to determine a node’s color directly from it’s column position and the global application parameters.

Optimizing locality through coloring

NABBITC requires that the user intentionally distribute data across their system and provides a coloring that captures computation locality. For this we make two assumptions about color: (1) that data initialized by each individual worker thread is given a unique color and (2) that each node of the computation task-graph is assigned a single color. Limiting nodes to a single color can lead to some information loss about a node’s locality. For example, a node can require data of different colors — in these scenarios, the user should specify the node’s color to be the one that maximizes locality for that node.

In most applications, a naive even distribution of data and a simple coloring function will, in fact, suffice. Consider the case of two dimensional stencils where the data is organized in a two dimensional array in row-major order. Data can be distributed by assigning an equal number of columns to each processor and colored accordingly. Figure 6 demonstrates how the subsequent task-coloring function would be defined.

IV. ANALYSIS OF COMPLETION TIME

We now present a simple analysis showing that the modifications made to NABBIT do not negatively effect the asymptotic runtime in a significant manner—this implies that NABBITC also provides almost asymptotically optimal load balancing for programs that have enough parallelism.

Just as in the NABBIT paper [1], say, we are given a task graph $G = (V, E)$, where each node u has work $W(u)$. Also, say that s is the unique node with zero in-degree and t is the unique node with zero out-degree. If these nodes are not unique, we can trivially add dummy root and final nodes. Define M as the number of nodes on the longest path in V from the source s to the sink t .

We can define the work T_1 as the time it takes to execute the task graph on a single processor and span T_∞ as the time it takes to execute it on an infinite number of processors. Therefore, the work is $T_1 = \sum_{u \in V} W(u) + O(|E|)$. The second term is due to the fact that each edge needs to be checked to make sure that it is satisfied. Similarly, the span

is $T_\infty = \max_{p \in \text{paths}(s,t)} \{\sum_{u \in p} W(u) + O(M)\}$ since nodes along any path through V can not execute in parallel. By the work and span laws [3, p. 780], the completion time on P processors for a task graph is at least $\max\{T_1/P, T_\infty\}$.

We will prove the following theorem—the analysis is a small extension to the analysis of runtime for NABBIT.

Theorem 1. *For task graph $G = (V, E)$ with maximum degree d , NABBITC executes G in time $O(T_1/P + T_\infty + M \lg d + \lg(P/\epsilon) + C)$ time on P processors with probability at least $1 - \epsilon$ where C is the amount of time each worker spends at startup trying to find a node of its own color.*

This theorem is similar to the theorem proved for NABBIT [1] apart from the last term C . The main difference between NABBIT and NABBITC is the fact of colored steals. In particular, when a worker runs out of work in NABBIT, it performs a random steal. On the other hand, when a worker runs out of work in NABBITC, it first checks a constant number of dequeues to see if it can find work of its own color and only performs a random steal if all these checks fail. In addition, at the start of the computation, NABBITC forces a colored steal and each processor may make C checks to find a node of its own color where C may be unbounded (in all our experiments, C turns out to be small.)

Lemma 1. *The total number of colored steals performed by NABBITC is $O(W+S+PC)$ where S is the number of random steal attempts, W is the number of steps the processors spend working on computation nodes, and C is the number of checks each processor performs at the beginning of the computation. Consequently, the total number of colored steals is bounded by $O(T_1 + PT_\infty + PM \lg d + P \lg(P/\epsilon) + PC)$*

Proof: Trivially, the number of checks at the beginning of the computation is PC since each processor performs at most C of them. After this, after a constant number of checks, a processor has either found work (therefore, these checks are bounded by $O(W)$) or the processor performs a random steal (these checks are bounded by $O(S)$). Summing these up gives us the result. The NABBIT analysis proves that the total number of work steps in the computation is at most T_1 and the total number of steal steps is at most $O(PT_\infty + PM \lg d + P \lg(P/\epsilon))$. This gives us the desired bound. ■

At any step, a worker is either working, doing a random steal, or doing a colored steal. Therefore, the total number of processor steps is bounded by $O(T_1 + PT_\infty + PM \lg d + P \lg(P/\epsilon) + PC)$ Since there are a total of P workers, we can divide by P to get the desired running time.

V. EXPERIMENTAL EVALUATION

We now compare NABBITC’s performance against original NABBIT and OPENMPTASKS and OPENMP loops. OPENMP offers multiple scheduling strategies for parallel for loops. The OPENMPSTATIC policy simply divides up the iteration space evenly among workers while OPENMPGUIDED dynamically load balances using adaptive block sizes. OPENMPTASKS is a tasking system in OpenMP which uses a centralized queue to load balance among workers. In particular, we try to answer the following questions:

- How well does NABBITC perform against NABBIT, OPENMP and OPENMPTASKS in general? We answer this question by using a set of compute and memory intensive applications, with either regular or irregular memory access patters. We find that the best scheduling strategy depends heavily on the application type— NABBITC shines best when both locality and load balance are necessary, while OPENMP edges out slightly over NABBITC when static schedules can get perfect locality and OPENMPTASKS is best in compute intensive applications where locality is not a determining factor for performance.
- To what extent does NABBITC improve data locality? We find that NABBIT has significantly fewer remote accesses compared to NABBIT and OPENMPTASKS while getting close to OPENMP when static scheduling can get perfect locality.
- Does the use of colored steals increase the overall cost to find work as compared to random stealing? We find that while the cost of enforcing the first colored steal is significant, NABBITC makes up for this overhead by having fewer steal attempts later.
- What is the impact of the choice of colors by the user? We consider the behavior of NABBITC using two particularly bad color choices and compare its behavior with NABBIT.

Experimental Setup: All our experiments were performed on an 80-core NUMA machine with 8 Intel Xeon E7-8860 2.27GHz 10-core processors and 1TB of collective DRAM. The machine uses Red Hat Linux 4.4.7-9 configured with 4KB pages. We use a stable GCC 4.9.0 build from the gcc-cilkplus branch for compiling our OPENMP and NABBIT benchmarks, extend this build for NABBITC and use a stable GCC 5.2.0 build for OPENMPTASKS⁴.

Benchmarks and Baselines: Table I details the benchmarks and input configurations used. The first five benchmarks exhibit regular memory access patterns. We consider these benchmarks to demonstrate the limitations of a dynamic task graph scheduler such as NABBIT (and OPENMPTASKS) that does not account for locality, and evaluate the potential for NABBITC to address these limitations. For these benchmarks, OPENMPSTATIC gets optimal locality (since the initialization and computation loops are perfectly matched) and good load balance since static load balancing is easy in these regular applications. Therefore, we only compare against this OPENMP strategy since it always performs better than OPENMPGUIDED.

The Smith-Waterman dynamic program [7] benchmarks exhibit highly regular memory access patterns and are more compute intensive than memory bound. We have implemented the wavefront computation in OPENMP, which must synchronize at each diagonal step. In NABBIT, NABBITC and OPENMPTASKS, we model the entire computation as a task-graph, exposing more parallelism.

PageRank iteratively computes the PageRank using the power method [8]. This benchmark exhibits access patterns dependent on the graph structure, with varying amounts of work per vertex. We consider three data sets from web crawls [4]

⁴OPENMPTASKS was added to OMPV4.0 which was included starting with GCC 5.0

Benchmark	Description	Problem size	Iterations	Task graph nodes	Serial	OPENMPSTATIC time (seconds)
cg	NAS conjugate gradient	$NA = 900000, NNZ = 26$	1	300		309
mg	NAS multigrid	$n\{x, y, z\} = 2048, LM = 11$	1	16384		690
heat	Heat diffusion stencil	$n = 16384, m = 655360$	5	102400		377
fdtd	Finite difference time domain	$n = 16384, m = 655360$	5	102400		970
life	Conway’s game of life	$n = 16384, m = 655360$	5	102400		275
sw	Smith-Waterman (n^3)	$\{n, m\} = 5120, B = 32x32$	1	25600		935
swn2	Smith-Waterman (n^2)	$\{n, m\} = 131072, B = 1024x1024$	1	16384		179
page-uk-2002	PageRank (power method) uk-2002 dataset	$nv = 18M, ne = 298M$	10	1800		198
page-uk-2007-05	PageRank (power method) uk-2007-05 dataset	$nv = 105M, ne = 3738M$	10	10500		960
page-twitter-2010	PageRank (power method) twitter-2010 dataset	$nv = 41M, ne = 1468M$	10	4100		1025

TABLE I: Benchmark configurations and serial OPENMPSTATIC execution time. The PageRank benchmarks use the same code with three different web crawl datasets [4], [5], [6].

that vary in size and graph structure. Specifically, twitter-2010 shows wider variation in its connectivity (e.g., much larger maximum out-degree) than the other data sets considered. On this benchmark, we compare against both OPENMPSTATIC and OPENMPGUIDED strategies.

Coloring strategy: In all benchmarks, we used OPENMP to distribute data evenly across the machine, with each processor core initializing a unique region of the data. Each thread is pinned to a processor core and given a unique color. During initialization, each data region is colored based on the color of the thread that initializes it. In order to color each node, we pick the color corresponding to the largest fraction of data as the node’s color. Computing the largest color is expensive for irregular benchmarks such as PageRank, where the accesses are data-dependent. In PageRank, each task takes a block of pages as input, which are accessed regularly, and updates the ranks of pages linked to them, which are accessed in an irregular fashion. The irregular accesses while traversing the links are not avoidable. Therefore, we color each task based on the block of pages it takes as input.

A. Overall performance

We now demonstrate the effect of locality-guided scheduling on the overall performance. In Figure 7, we present the speedup achieved by OPENMP, OPENMPTASKS, NABBIT, and NABBITC over serial execution. Error bars show standard deviation across five runs.

In the first 5 (regular) applications, NABBITC performs better than NABBIT and OPENMPTASKS due to getting better data locality. We see that in cg, when there are very few nodes in the task graph, NABBITC’s benefit over original NABBIT becomes negligible because processor cores have few nodes to work with. With mg, heat, fdtd, and life, when there are many nodes in the task graph, NABBITC is able to continue getting good performance while original NABBIT and OPENMPTASKS suffer due to its locality-obliviousness. For these benchmarks, we see that OPENMPSTATIC consistently performs best (but NABBITC performance approaches it). When threads are pinned and the computation loops are scheduled in the same way as the data initialization loops, OPENMP achieves the maximum locality possible despite not having received any explicit locality hints from the programmer. In addition, it also achieves good load balance, since each iteration does approximately equal amount of work.

Smith-Waterman programs are compute intensive and load balance is more important than locality. Here, NABBITC, NABBIT and OPENMPTASKS all perform better than OPENMPSTATIC since they get better load balance and exploit more parallelism. Here, OPENMPTASKS performs slightly better than NABBIT and NABBITC due to lower scheduling overheads — all the tasks in the graph are relatively large leading to efficient centralized scheduling.

PageRank is an irregular application where both load balance and locality are important and here is where NABBITC shines. OPENMPSTATIC is not able to maintain good locality and load balance automatically. On the other hand, OPENMPSTATIC and OPENMPGUIDED do not perform well due to high scheduling overheads on fine-grained irregular tasks. NABBIT and NABBITC both perform comparably, but NABBITC is slightly better on the larger benchmarks due to better locality.

In summary, NABBITC almost always provides the best or the second best performance by providing a good trade-off between locality, load balance and low overhead. All the other systems perform better on some benchmarks but much worse on others since they either do static scheduling (which provides bad load balance on irregular applications) or only consider load balance.

B. Locality impact of NABBITC’s scheduling strategy

We now look closer at the locality achieved by NABBITC during the execution of these benchmarks. Because counting each memory reference might be expensive⁵, we perform this check at the node level in the task graph. This consists of two parts. Note that each of our evaluation system consists of eight NUMA domains, each with 10 cores. First, for each thread, we count the number of nodes it executes that are not the same color as any thread in the same NUMA node. Second, for each thread, we check all predecessors of executed nodes, and count those that are not the same color as any thread in the same NUMA node. Sum of these counts across all threads is reported as the number of remote accesses. For the regular benchmarks, we can compute this as the benchmarks execute without perturbing the execution. For PageRank, this instrumentation can significantly perturb the execution time. Therefore, we track the nodes executed by each thread to

⁵We were limited by OS version and available hardware counters and were unable to measure remote accesses, stall cycles, etc.

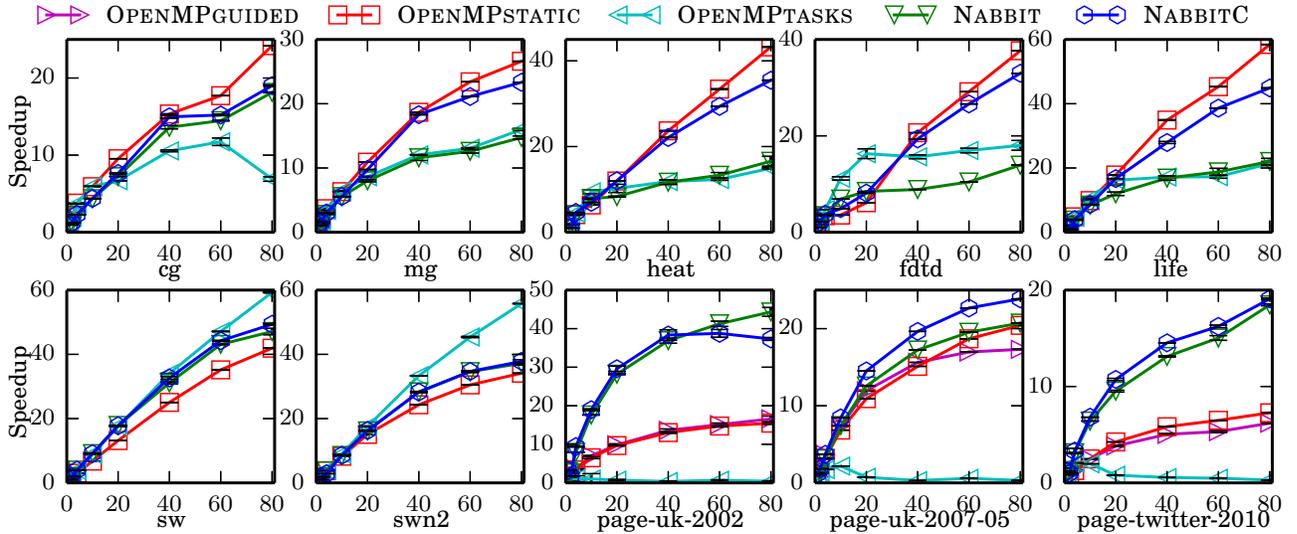


Fig. 7: Speedup for all benchmarks. x-axis: number of threads (processor cores); y-axis: speedup over serial OPENMPSTATIC.

record the schedule used in the timing runs. This schedule is replayed to compute the percentage of remote accesses.

Figure 8 shows the percentage of accesses that are remote for NABBIT, NABBITC, OPENMPSTATIC, and OPENMPTASKS on 20 or more processor cores (smaller core counts fit in one NUMA domain and do not incur remote accesses). Since both OPENMPTASKS and NABBIT do not consider locality while scheduling, they incur comparable and significant number (44–88%) of remote accesses. The introduction of colored steals in NABBITC significantly decreases the percentage of remote accesses. For all benchmarks except twitter-2010 and the Smith-Waterman benchmarks, NABBITC incurs 0% to 9% remote accesses. Importantly, unlike in the case of NABBIT, this percentage does not strictly increase with the number of processors. All strategies incur a high percentage of remote accesses for twitter-2010 and Smith-Waterman. The twitter-2010 dataset has a particularly large out-degree with many edges connecting distant nodes, resulting in a large number of unavoidable remote accesses. The Smith-Waterman applications suffer from a similar problem, where unavoidable accesses are baked into the algorithm.

For regular applications, OPENMPSTATIC incurs almost no remote accesses, as we expect from how the data is initialized. For PageRank, OPENMPSTATIC still has fewer remote accesses than NABBITC; however, as we saw above, it does not have good performance since it is unable to provide good load balance. This result indicates the importance of both locality and load balance—while NABBIT provides great load balance and OPENMPSTATIC provides great locality, NABBITC performs better than both on this irregular benchmark since it simultaneously considers both metrics.

C. Overheads due to colored steals

The two sources of overhead for NABBITC arise from requiring a constant number of colored steals before performing random steals and forcing the first steal to be a colored steal.

Effect on total steals: We now look at the comparison of NABBITC and NABBIT at a more fine-grained level. In Figure 9 we see that NABBITC, perhaps counter-intuitively, performs far fewer total successful steals than NABBIT. The introduction of colored steals, and specifically enforcing the first colored steal, helps to significantly reduce the total number steals by ensuring that thieves acquire nodes higher up in the task graph to start with. Due to the depth-first nature of the scheduler, nodes higher up in the task graph have more potential work. Therefore, by ensuring thieves begin with nodes connected to the root, NABBITC is able to effectively increase the amount of work each worker begins with, reducing the total number of steals required.

Overhead due to enforcing first colored steal: To calculate the overhead of ensuring that the first steal is a colored steal, Fig. 10 shows the average amount of time processor cores spent waiting to acquire work for the heat benchmark. We observed that the times were very similar for all other benchmarks and do not present them here due to space limitations. While this overhead can be substantial, it is agnostic to the application, provided there is at least one node from each color connected to the root. This startup cost can be amortized out with larger, longer running benchmarks. Additionally, recall that we observed that, in practice, enforcing the first colored steal results in far fewer total number of steals which makes up for this overhead.

D. Importance of good coloring

Overheads with invalid coloring: To evaluate this worst case overhead from attempted colored steals, we assigned all nodes an invalid color (no worker has this color), ensuring that all colored steals fail. Therefore, this version of NABBITC behaves like original NABBIT apart from incurring the overheads of colored steals. In Table. III, we see that NABBITC with this alternative coloring performs comparable to original NABBIT. This leads to two conclusions. First, NABBITC achieves all its performance advantage over NABBIT via coloring and not

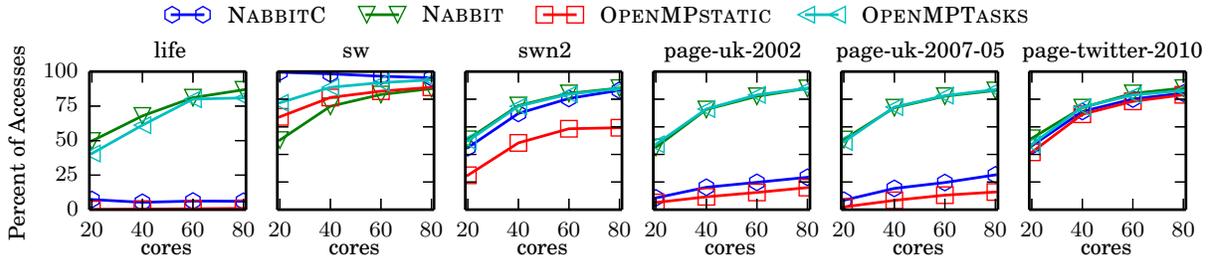


Fig. 8: Percentage of accesses that are data in remote NUMA domains. We show percentages for 20–80 cores (1–10 cores fit in one NUMA domain and do not incur remote accesses). We’ve omitted cg, mg, heat and ftdt as their trends are identical to that of life

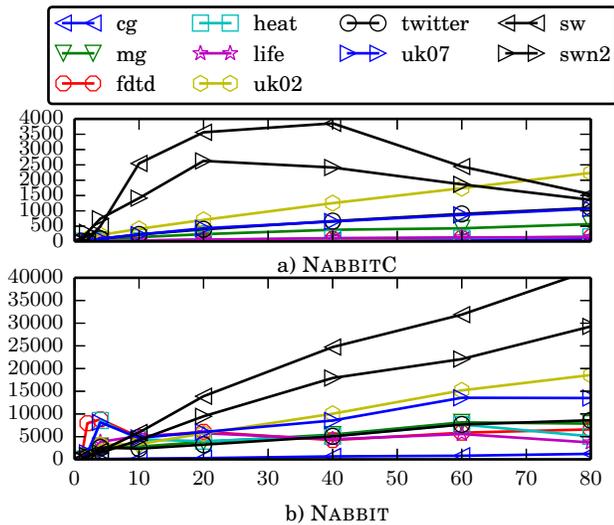


Fig. 9: Average number of successful steals for (a) NABBITC and (b) NABBIT. x-axis: number of cores; y-axis: Average number of successful steals

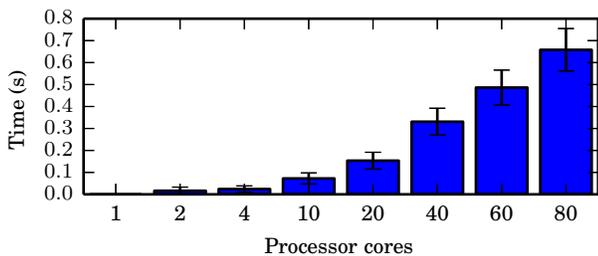


Fig. 10: Average idle time per processor core (across all processor cores and runs) due to forcing the first colored steal for the heat benchmark. Error bars show standard deviation across five runs among all processor cores. We observed this time was the same for all benchmarks.

via some accidental optimization. Second, the additional work performed by colored steals introduces minimal overhead — the mean speedups are within one or two standard deviations, indicating that, for the benchmarks considered, colored steals incur no statistically significant overhead.

P	cg	mg	heat	ftdt	life	uk-02	twitter	uk-07	sw	swn2
20	0.96	0.76	0.66	1.08	0.65	0.75	0.83	0.76	0.76	0.66
40	0.97	0.88	0.78	1.10	0.77	0.83	0.55	0.94	0.88	0.78
60	0.96	0.97	0.95	1.08	0.94	0.89	0.49	1.02	0.97	0.95
80	1.02	0.96	0.93	0.94	0.97	0.91	0.41	1.02	0.96	0.93
	0.98	0.89	0.83	1.05	0.83	0.85	0.57	0.93	0.89	0.83
	0.02	0.08	0.12	0.06	0.13	0.06	0.16	0.11	0.08	0.12

TABLE II: Speedup of NABBITC over NABBIT when all tasks are assigned bad colors resulting in preferential execution of non-local tasks.

P	cg	mg	heat	ftdt	life	uk02	twitter	uk07	sw	swn2
20	1.03	0.99	0.94	1.06	0.93	1.03	1.12	1.12	0.99	0.94
40	1.02	0.99	0.99	1.04	0.92	0.97	1.09	1.10	0.99	0.99
60	0.99	0.98	0.94	1.03	0.92	0.98	1.01	1.08	0.98	0.94
80	1.06	0.97	0.88	0.91	0.98	0.94	1.07	1.07	0.97	0.88
	1.03	0.99	0.94	1.01	0.94	0.98	1.07	1.09	0.99	0.94
	0.02	0.01	0.04	0.06	0.03	0.03	0.04	0.02	0.01	0.04

TABLE III: Speedup of NABBITC over NABBIT when all tasks are assigned invalid colors resulting in failure of all colored steal attempts.

Behavior under bad coloring: The performance of NABBITC is directly tied to the coloring provided by the user. NABBITC assumes the user has constructed a “good” coloring and makes decisions based on this assumption. We now test the importance of good coloring in another way — we force all workers to preferentially perform non-local work by creating a coloring where all nodes are given valid incorrect colors. The data in Table. II shows that NABBITC with this bad coloring performs similar to NABBIT indicating that a good coloring is important. More interestingly, it shows that NABBIT provides really bad locality since NABBITC with intentionally bad locality performs within two standard deviations of NABBIT.

VI. RELATED WORK

Static task graph schedulers [9], [10], [11] minimize completion time while maximizing locality [12] by completing expanding and analyzing a task graph, together with accurate information on computation and communication costs associated with each task. We consider task graphs that are

dynamically explored and do not require prior knowledge of task and communication times.

Xkaapi [?] introduces a number heuristics to improve OpenMP task scheduling. A work-stealing scheduler utilizes user specified data distribution patterns, task dependences and a model of the machine topology to guide scheduling decisions. Task graphs are explored upwards, with workers pushing successor tasks onto their local work queue. These tasks are not immediately ready and thieves are required to scan a victim's work queue and resolve dependences in order to find ready tasks.

Cilk's random work stealing is agnostic of locality considerations [13]. Several efforts have incorporated locality considerations by altering the work stealing strategy [14], [15], [16], [17], [18]. These approaches do not naturally extend to scheduling data-flow graphs while preserving provably efficiency in terms of scheduling overheads and effectiveness of load balancing. Locality-aware work stealing strategies (e.g., SLAW [14]) focus on improving stack locality rather than locality of data in heap, the focus of this paper.

Event-driven scheduling strategies map tasks to locality domains together with efficient identification and tracking of ready tasks that can be scheduled [19], [20]. In these systems, data distribution implies a computation partitioning with no further migration of tasks to tackle load imbalance.

SuperMatrix [21], a runtime scheduling system for algorithms operating on blocks as observed in linear algebra programs, mimics a superscalar microarchitecture's scheduling strategy in software. It does not account for data locality. Dague [22], a distributed DAG engine, improves locality by working on the local queue when possible. XKaapi [23] is a work-stealing-based scheduler for task graphs that pushes tasks to processors that have better locality for those tasks. StarPU [24], a task-graph scheduler for heterogeneous multi-core systems, incorporates multiple locality-aware schedulers. These approaches, while improving data locality, do not preserve the critical path length or provide provable parallel efficiency.

SMPSs [2] schedules dependent tasks together to improve locality. Legion [25] exploits user-specified locality information and coherence properties to perform locality-aware scheduling using a software out-of-order processor. CnC [26] allows the specification of task graphs that are scheduled using a variety of strategies. Legion and CnC also allow user-specification to control task mapping and scheduling (using mappers in Legion and tuners in CnC). Olivier et al. developed various strategies to schedule OpenMP tasks including hierarchical scheduling, and work stealing by one thread on behalf of others in the same chip [27]. None of these schedulers in these systems attempt to preserve optimality guarantees. However, the scheduling strategy developed in this paper can be used to develop provably efficient and locality-aware scheduling algorithms for these task-graph frameworks.

Bugnion et al. [28] developed compiler-directed page coloring techniques to minimize conflict misses. Chilimbi and Shaham [29] identified *hot* data streams and then colocated them to improve spatial locality. Chen et al. studied scheduling threads for constructive cache sharing [30]. Various ap-

proaches have studied the partitioning of shared caches among threads (e.g., [31], [32], [33]). These approaches cannot be applied to optimize NUMA locality considered in this paper.

VII. CONCLUSIONS AND FUTURE WORK

In this paper, we have presented NABBITC, a flexible and easy-to-use task graph library that allows the user to provide locality hints via the use of coloring and provides good load balance via dynamic scheduling. NABBITC is geared towards scheduling on NUMA hardware, where remote accesses may be considerably more expensive than local accesses, but one must strike a balance between locality and load balance to get good performance. Experimental results indicate that this approach is promising, especially for memory intensive irregular applications running on NUMA machines, where static scheduling can compromise load balancing and locality-unaware dynamic scheduling has too many remote accesses. While NABBITC uses Cilk Plus as the underlying language and runtime, we believe this approach can be implemented on other systems such as Intel's Threading Building Blocks.

For future work, we would like to explore automatically discovering colors, reducing the user's workload and thus providing these locality benefits for free. NABBITC's startup overhead could also be significantly reduced with additional assistance from the underlying runtime system which could allow it to even remain competitive with OPENMPTASKS when locality is less important.

ACKNOWLEDGEMENTS

This research was supported by National Science Foundation grants CCF-1527692, CCF-143906, and CCF-1150036, and by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research under award number 63823. Pacific Northwest National Laboratory is operated by Battelle for DOE under Contract DE-AC05-76RL01830.

REFERENCES

- [1] K. Agrawal, C. E. Leiserson, and J. Sukha, "Executing task graphs using work-stealing," in *IPDPS*, 2010, pp. 1–12.
- [2] J. M. Perez, R. M. Badia, and J. Labarta, "A dependency-aware task-based programming environment for multi-core architectures," in *CLUSTER*, 2008, pp. 142–151.
- [3] T. Cormen, C. Leiserson, R. Rivest, and C. Stein, *Introduction to Algorithms*. MIT Press, 2009.
- [4] "Laboratory for web algorithmics," <http://law.di.unimi.it/datasets.php>.
- [5] P. Boldi and S. Vigna, "The WebGraph framework I: Compression techniques," in *WWW*, 2004, pp. 595–601.
- [6] P. Boldi, M. Rosa, M. Santini, and S. Vigna, "Layered label propagation: A multiresolution coordinate-free ordering for compressing social networks," in *WWW*, 2011, pp. 587–596.
- [7] T. F. Smith and M. S. Waterman, "Identification of common molecular subsequences," *Journal of molecular biology*, vol. 147, no. 1, pp. 195–197, 1981.
- [8] L. Page, S. Brin, R. Motwani, and T. Winograd, "The PageRank citation ranking: bringing order to the web," 1999.
- [9] Y. Kwok and I. Ahmad, "Static scheduling algorithms for allocating directed task graphs to multiprocessors," *ACM Comput. Surv.*, vol. 31, no. 4, pp. 406–471, Dec. 1999.
- [10] G. Liu, K.-L. Poh, and M. Xie, "Iterative list scheduling for heterogeneous computing," *JPDC*, vol. 65, no. 5, pp. 654–665, 2005.
- [11] T. Yang and A. Gerasoulis, "PYRROS: static task scheduling and code generation for message passing multiprocessors," in *SC*, 1992, pp. 428–437.

- [12] N. Vydyanathan, S. Krishnamoorthy, G. M. Sabin, Ü. V. Çatalyürek, T. M. Kurç, P. Sadayappan, and J. H. Saltz, "An integrated approach to locality-conscious processor allocation and scheduling of mixed-parallel applications," *TPDS*, vol. 20, no. 8, pp. 1158–1172, 2009.
- [13] M. Frigo, C. Leiserson, and K. Randall, "The implementation of the Cilk-5 multithreaded language," in *PLDI*, 1998, pp. 212–223.
- [14] Y. Guo, J. Zhao, V. Cave, and V. Sarkar, "SLAW: a scalable locality-aware adaptive work-stealing scheduler for multi-core systems," in *IPDPS*, 2010, pp. 341–342.
- [15] S.-J. Min, C. Iancu, and K. Yelick, "Hierarchical work stealing on manycore clusters," in *PGAS*, 2011.
- [16] J.-N. Quintin and F. Wagner, "Hierarchical work-stealing," in *Euro-Par 2010-Parallel Processing*, 2010, pp. 217–229.
- [17] J. Lifflander, S. Krishnamoorthy, and L. V. Kale, "Optimizing data locality for fork/join programs using constrained work stealing," in *SC*, 2014, pp. 857–868.
- [18] J. Lifflander and S. Krishnamoorthy, "Cache locality optimization for recursive programs," in *PLDI*, 2017.
- [19] S. Chatterjee, S. Tasirlar, Z. Budimlic, V. Cavé, M. Chabbi, M. Grossman, V. Sarkar, and Y. Yan, "Integrating asynchronous task parallelism with MPI," in *IPDPS*, 2013, pp. 712–725.
- [20] G. Bosilca, A. Bouteiller, A. Danalis, M. Faverge, T. Héroult, and J. J. Dongarra, "PaRSEC: Exploiting heterogeneity to enhance scalability," *Computing in Science & Engineering*, vol. 15, no. 6, pp. 36–45, 2013.
- [21] E. Chan, E. S. Quintana-Orti, G. Quintana-Orti, and R. Van De Geijn, "SuperMatrix out-of-order scheduling of matrix operations for smp and multi-core architectures," in *SPAA*, 2007, pp. 116–125.
- [22] G. Bosilca, A. Bouteiller, A. Danalis, T. Héroult, P. Lemarinier, and J. Dongarra, "Dague: A generic distributed dag engine for high performance computing," *ParCo*, vol. 38, no. 1, pp. 37–51, 2012.
- [23] T. Gautier, J. V. Lima, N. Maillard, and B. Raffin, "XKaapi: A runtime system for data-flow task programming on heterogeneous architectures," in *IPDPS*, 2013, pp. 1299–1308.
- [24] C. Augonnet, S. Thibault, R. Namyst, and P.-A. Wacrenier, "StarPU: a unified platform for task scheduling on heterogeneous multicore architectures," *CCPE*, vol. 23, no. 2, pp. 187–198, 2011.
- [25] M. Bauer, S. Treichler, E. Slaughter, and A. Aiken, "Legion: expressing locality and independence with logical regions," in *SC*, 2012, p. 66.
- [26] Z. Budimlić, M. Burke, V. Cavé, K. Knobe *et al.*, "Concurrent collections," *Scientific Programming*, vol. 18, no. 3-4, pp. 203–217, 2010.
- [27] S. L. Olivier, A. K. Porterfield, K. B. Wheeler, and J. F. Prins, "Scheduling task parallelism on multi-socket multicore systems," in *ROSS*, 2011, pp. 49–56.
- [28] E. Bugnion, J. M. Anderson, T. C. Mowry, M. Rosenblum, and M. S. Lam, "Compiler-directed page coloring for multiprocessors," in *ACM SIGPLAN Notices*, vol. 31, no. 9. ACM, 1996, pp. 244–255.
- [29] T. M. Chilimbi and R. Shaham, "Cache-conscious coallocation of hot data streams," in *PLDI*, 2006, pp. 252–262.
- [30] S. Chen, P. B. Gibbons, M. Kozuch, V. Liaskovitis, A. Ailamaki, G. E. Blelloch *et al.*, "Scheduling threads for constructive cache sharing on cmps," in *SPAA*, 2007, pp. 105–115.
- [31] M. K. Qureshi and Y. N. Patt, "Utility-based cache partitioning: A low-overhead, high-performance, runtime mechanism to partition shared caches," in *MICRO*, 2006, pp. 423–432.
- [32] G. Suh, L. Rudolph, and S. Devadas, "Dynamic partitioning of shared cache memory," *Journal of Supercomp.*, vol. 28, no. 1, pp. 7–26, 2004.
- [33] F. Mueller, "Compiler support for software-based cache partitioning," in *ACM SIGPLAN Notices*, vol. 30, no. 11. ACM, 1995, pp. 125–133.