

# A scalable clustering-based task scheduler for homogeneous processors using DAG partitioning

Yusuf Özkaya, Anne Benoit, Bora Uçar, Julien Herrmann, Umit Catalyurek

## ► To cite this version:

Yusuf Özkaya, Anne Benoit, Bora Uçar, Julien Herrmann, Umit Catalyurek. A scalable clustering-based task scheduler for homogeneous processors using DAG partitioning. IPDPS 2019 - 33rd IEEE International Parallel & Distributed Processing Symposium, May 2019, Rio de Janeiro, Brazil. pp.155-165, 10.1109/IPDPS.2019.00026 . hal-02082794

**HAL Id: hal-02082794**

**<https://hal.inria.fr/hal-02082794>**

Submitted on 28 Mar 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A scalable clustering-based task scheduler for homogeneous processors using DAG partitioning

M. Yusuf Özkaya\*, Anne Benoit\*<sup>†</sup>, Bora Uçar\*<sup>†</sup>, Julien Herrmann\*, and Ümit V. Çatalyürek\*

\*Georgia Institute of Technology, School of Computational Science and Engineering, Atlanta, GA, USA

<sup>†</sup>Univ. Lyon, CNRS, ENS de Lyon, Inria, Univ. Claude-Bernard Lyon 1, LIP UMR5668, France

Email: {myozka, jherrmann9, umit}@gatech.edu, {anne.benoit, bora.ucar}@ens-lyon.fr

**Abstract**—When scheduling a directed acyclic graph (DAG) of tasks with communication costs on computational platforms, a good trade-off between load balance and data locality is necessary. List-based scheduling techniques are commonly-used greedy approaches for this problem. The downside of list-scheduling heuristics is that they are incapable of making short-term sacrifices for the global efficiency of the schedule. In this work, we describe new list-based scheduling heuristics based on clustering for homogeneous platforms, under the realistic duplex single-port communication model. Our approach uses an acyclic partitioner for DAGs for clustering. The clustering enhances the data locality of the scheduler with a global view of the graph. Furthermore, since the partition is acyclic, we can schedule each part completely once its input tasks are ready to be executed. We present an extensive experimental evaluation showing the trade-offs between the granularity of clustering and the parallelism, and how this affects the scheduling. Furthermore, we compare our heuristics to the best state-of-the-art list-scheduling and clustering heuristics, and obtain more than three times better makespan in cases with many communications.

**Index Terms**—List scheduling, clustering, partitioning, directed acyclic graphs, data locality, concurrency.

## I. INTRODUCTION

Scheduling is one of the most studied areas of computer science. A large body of research deals with scheduling applications/workflows modeled as Directed Acyclic Graphs (DAGs), where vertices represent atomic tasks, and edges represent dependencies with associated communication costs [17], [24]. The classical objective function is to minimize the total execution time, or *makespan*, and this problem is denoted as  $P|prec, c_{i,j}|C_{max}$  in the scheduling literature. Among others, *list-based scheduling* techniques are the most widely studied and used techniques, mainly due to the ease of implementation and explanation of the progression of the heuristics [1], [12], [18], [21], [23], [26], [27], [28]. In list-based scheduling techniques, tasks are ordered based on some predetermined priority, and then are mapped and scheduled onto processors. Another widely used approach is clustering-based scheduling [14], [15], [22], [27], [28], [29], where tasks are grouped into clusters and then scheduled onto processors.

Almost all of the existing clustering-based scheduling techniques are based on bottom-up clustering approaches, where clusters are constructively built from the composition of atomic tasks and existing clusters. We argue that such decisions are local, and hence cannot take into account the global structure of the graph. Recently, we have developed one of the

first multi-level acyclic DAG partitioners [11]. We hypothesize that clusters found using such a DAG partitioner are much more successful in putting together the tasks with complex dependencies, and hence in minimizing the overall inter-processor communication, and we confirm this hypothesis in our experiments.

In this work, we use the realistic duplex single-port communication model, where at any point in time, each processor can, in parallel, execute a task, send one data, and receive another data. Because concurrent communications are limited within a processor, minimizing the communication volume is crucial to minimizing the total execution time (makespan).

The goal is to minimize the makespan when the DAG is executed on a parallel platform. In our proposed schedulers, when scheduling to a system with  $p$  processing units (or processors), the original task graph is first partitioned into  $K$  parts (clusters), where  $K \geq p$ . Then, a list-based scheduler is used to assign tasks (not the clusters). Our scheduler hence uses list-based scheduler, but with one major constraint: all the tasks of a cluster will be executed by same processor. This is not the same as scheduling the graph of clusters, as the decision to schedule a task can be made before scheduling all tasks in a predecessor cluster. Our intuition is that, since the partition is done beforehand, the scheduler “sees” the global structure of the graph, and it uses this to “guide” the scheduling decisions. Since all the tasks in a cluster will be executed on the same processor, the execution time for the cluster can be approximated by simply the sum of the individual tasks’ weights (actual execution time can be larger due to dependencies to tasks that might be assigned to other processors). Here, we heuristically decide that having balanced clusters helps the scheduler to achieve load-balanced execution. The choice of the number of parts  $K$  is a trade-off between data locality vs. concurrency. Large  $K$  values may yield higher concurrency, but would potentially incur more inter-processor communication. At the extreme, each task is a cluster, where we have the maximum potential concurrency. However, in this case, one has to rely on list-based scheduler’s local decisions to improve data locality, and hence reduce inter-processor communication.

Our main contribution is to develop three different variants (meta-heuristics) of partitioning-assisted list-based scheduler, taking different decisions about how to schedule tasks within a part. These variants run on top of two classical list-based

schedulers: (1) BL-EST chooses the task with largest bottom-level first (BL), and assigns the task on the processor with the earliest start time (EST), while (2) ETF tries all ready tasks on all processors and picks the combination with the earliest EST first (hence with a higher complexity). The proposed meta-heuristics can be used with any other list scheduler and DAG partitioner, hence they provide a flexible solution to DAG scheduling. Also, we experimentally evaluate the new algorithms against the two baseline list-based schedulers (BL-EST and ETF) and one baseline cluster-based scheduler (DSC-GLB-ETF), since ETF and DSC-GLB-ETF are the winners of the recent comparison done by Wang and Sinnen [27]. However, unlike [27], we follow the realistic duplex single-port communication model. We show significant savings in terms of makespan, in particular when the communication-to-computation ratio (CCR) is large, i.e., when communications matter a lot, hence demonstrating the need for a partitioning-assisted scheduling technique.

In other words, we propose a novel algorithmic framework for DAG scheduling, building upon a multi-level *acyclic* DAG partitioner for the clustering phase. Furthermore, we consider a realistic communication model, contrarily to most theoretical work on scheduling. Thus, our algorithms lend themselves as efficient heuristics with no lower bounds or performance guarantees. However, as demonstrated in the results section, they drastically outperform state-of-the-art schedulers under more realistic scenarios, such as single-port communication model and when communications are more costly than computations. For example, one of the datasets we experimented includes several DAGs corresponding to high-performance computing (HPC) applications that use Open Community Runtime (OCR) framework [30], on which we achieve more than three times better makespan than the state-of-the-art heuristic with large CCRs.

The rest of the paper is organized as follows. First, we discuss related work in Section II. Next, we introduce the model and formalize the optimization problem in Section III. The proposed scheduling heuristics are described in Section IV, and they are evaluated through extensive simulations in Section V. Finally, we conclude and give directions for future work in Section VI.

## II. RELATED WORK

Task graph scheduling has been the subject of a wide literature, ranging from theoretical studies to practical ones. Kwok and Ahmad [17] give an excellent survey and taxonomy of task scheduling methods and some benchmarking techniques to compare these methods [16].

On the theoretical side, a related problem of minimizing the makespan of a DAG on identical processors without communication costs ( $P|prec|C_{max}$ ) has been extensively studied. Graham's seminal list-scheduling algorithm [9] has been known for a long time to be a  $(2 - \frac{1}{p})$ -approximation algorithm, where  $p$  is the number of processors. It has then been shown that it is NP-hard to improve upon this approximation ratio, assuming a new variant of the unique games

conjecture [25]. Several works further focus on unit execution times to derive theoretical results (lower bounds, complexity results), see for instance [13].

On the practical side, communication costs cannot be neglected, and it becomes much harder to derive theoretical guarantees. Even the problem with unit execution time and unit communication time (UET-UCT) is NP-hard [20]. Hence, the  $P|prec, c_{i,j}|C_{max}$  problem is usually tackled through heuristics. For coarse-grain graphs, a guaranteed heuristic based on a linear programming formulation of the problem was proposed [10], and it was proven that there always exists a linear optimal clustering [7].

DAG scheduling heuristics can be divided into two groups with respect to whether they allow task duplication or not [2]. Those that allow task duplication do so to avoid communication. The focus of this work is non-duplication based scheduling. There are two main approaches taken by the non-duplication based heuristics: list scheduling and cluster-based scheduling. A recent comparative study [27] gives a catalog of list-scheduling and cluster-scheduling heuristics and compares their performance. These algorithms take the entire task graph as input, similar to our approach.

In the list-based scheduling approach [1], [9], [12], [18], [21], [23], [26], [28], each task in the DAG is first assigned a priority. Then, the tasks are sorted in descending order of priorities, thereby resulting in a priority list. Finally, the tasks are scheduled in topological order, with the highest priorities first. There are also two variants of list-scheduling based on how priorities are computed: *static* and *dynamic*. In the static list-scheduling, priorities are pre-computed and do not change during the algorithm. In the dynamic list-scheduling, task priorities are updated as the predecessor tasks are scheduled. The list-scheduling based heuristics usually are easy to implement and understand. In general, the static list-scheduling algorithms also have low computational complexity, whereas dynamic list-scheduling algorithms have higher complexity, due to priority updates.

In the cluster-based scheduling approach [14], [15], [22], [27], [28], [29], the tasks are first divided into clusters, each to be scheduled on the same processor. The clusters usually consist of highly communicating tasks. Then, the clusters are scheduled onto an unlimited number of processors, which are finally combined to yield the available number of processors.

Our approach is close to cluster-based scheduling in the sense that we first partition tasks into  $K \geq p$  clusters, where  $p$  is the number of available processors. At this step, we enforce somewhat balanced clusters. In the next step, we schedule tasks as in the list-scheduling approach, not the clusters, since there is a degree of freedom in scheduling a task of a cluster. Hence, our approach can also be conceived as a hybrid list and cluster scheduling, where the decisions of the list-scheduling part are constrained by the cluster-scheduling decisions.

We consider homogeneous computing platforms, where the processing units are identical and communicate through a homogeneous network. Task graphs and scheduling approaches can also be used to model and execute workflows on grids

and heterogeneous platforms [5], [8]; HEFT (heterogeneous earliest finish time) [26] is a common approach for this purpose. Assessing the performance of our new scheduling strategies on heterogeneous platforms will be considered in future work.

### III. MODEL

Let  $G = (V, E)$  be a directed acyclic graph (DAG), where the vertices in the set  $V$  represent tasks, and the edges in the set  $E$  represent the precedence constraints between those tasks. Let  $n = |V|$  be the total number of tasks. We use  $\text{Pred}[v_i] = \{v_j \mid (v_j, v_i) \in E\}$  to represent the (immediate) predecessors of a vertex  $v_i \in V$ , and  $\text{Succ}[v_i] = \{v_j \mid (v_i, v_j) \in E\}$  to represent the (immediate) successors of  $v_i$  in  $G$ . Vertices without any predecessors are called *source* nodes, and the ones without any successors are called *target* nodes. Every vertex  $v_i \in V$  has a weight, denoted by  $w_i$ , and every edge  $(v_i, v_j) \in E$  has a cost, denoted by  $c_{i,j}$ .

The computing platform is a homogeneous cluster consisting of  $p$  identical processing units, called *processors*, and denoted  $P_1, \dots, P_p$ , communicating through a fully-connected homogenous network. Each task needs to be scheduled onto a processor respecting the precedence constraints, and tasks are non-preemptive and atomic: a processor executes a single task at a time. For a given mapping of the tasks onto the computing platform, let  $\mu(i)$  be the index of the processor on which task  $v_i$  is mapped, i.e.,  $v_i$  is executed on the processor  $P_{\mu(i)}$ . For every vertex  $v_i \in V$ , its weight  $w_i$  represents the time required to execute the task  $v_i$  on any processor. Furthermore, if there is a precedence constraint between two tasks mapped onto two different processors, i.e.,  $(v_i, v_j) \in E$  and  $\mu(i) \neq \mu(j)$ , then some data must be sent from  $P_{\mu(i)}$  to  $P_{\mu(j)}$ , and this takes a time represented by the edge cost  $c_{i,j}$ .

We enforce the realistic duplex single-port communication model, where at any point in time, each processor can, in parallel, execute a task, send one data, and receive another data. Consider the DAG example in Figure 1, where all execution times are unitary, and communication times are depicted on the edges. The computing platform in the example of Figure 1 has two identical processors. There is no communication cost to pay when two tasks are executed on the same processor, since the output can be directly accessed in the processor memory by the next task. For the proposed schedule,  $P_1$  is already performing a *send* operation when  $v_5$  would like to initiate a communication, and hence this communication is delayed by 0.5 time unit, since it can start only after  $P_1$  has completed the previous send from  $v_1$  to  $v_2$ . However,  $P_1$  can receive data from  $v_2$  to  $v_3$  in parallel to sending data from  $v_5$  to  $v_6$ . In this example, the total execution time, or *makespan*, is 6.

Formally, a schedule of graph  $G$  consists of an assignment of tasks to processors (already defined as  $\mu(i)$ , for  $1 \leq i \leq n$ ), and a start time for each task,  $\text{st}(i)$ , for  $1 \leq i \leq n$ . Furthermore, for each precedence constraint  $(v_i, v_j) \in E$  such that  $\mu(i) \neq \mu(j)$ , we must specify the start time of the communication,  $\text{com}(i, j)$ . Several constraints must be

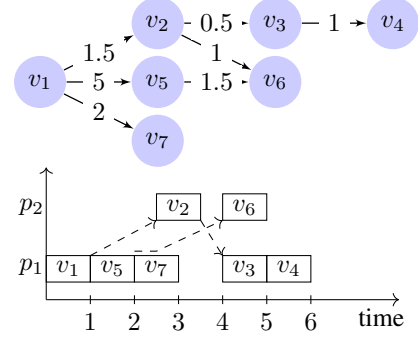


Figure 1: Example of a small DAG with seven vertices executed on a homogeneous platform with two processors.

Notation	Meaning
DAG	Directed Acyclic Graph
CCR	Communication-to-Computation Ratio
BL	Bottom-Level
TL	Top-Level
EST	Earliest Start Time
ETF	Earliest EST First
DSC	Dominant Sequence Clustering
GLB	Guided Load Balancing

Table I: Acronyms.

met to have a valid schedule, in particular with respect to communications:

- (atomicity) For each processor  $P_k$ , for each task  $v_i$  such that  $\mu(i) = k$ , the intervals  $[\text{st}(i), \text{st}(i) + w_i]$  are disjoint.
- (precedence constraints, same processor) For each  $(v_i, v_j) \in E$  with  $\mu(i) = \mu(j)$ ,  $\text{st}(i) + w_i \leq \text{st}(j)$ .
- (precedence constraints, different processors) For each  $(v_i, v_j) \in E$  with  $\mu(i) \neq \mu(j)$ ,  $\text{st}(i) + w_i \leq \text{com}(i, j)$  and  $\text{com}(i, j) + c_{i,j} \leq \text{st}(j)$ .
- (one-port, sending) For each  $P_k$ , for each  $(v_i, v_j) \in E$  such that  $\mu(i) = k$  and  $\mu(j) \neq k$ , the intervals  $[\text{com}(i, j), \text{com}(i, j) + c_{i,j}]$  are disjoint.
- (one-port, receiving) For each  $P_k$ , for each  $(v_i, v_j) \in E$  such that  $\mu(i) \neq k$  and  $\mu(j) = k$ , the intervals  $[\text{com}(i, j), \text{com}(i, j) + c_{i,j}]$  are disjoint.

The goal is then to minimize the makespan, that is the maximum execution time:

$$M = \max_{1 \leq i \leq n} \{\text{st}(i) + w_i\}. \quad (1)$$

We are now ready to formalize the MINMAKESPAN optimization problem: *Given a weighted DAG  $G = (V, E)$  and  $p$  identical processors, the MINMAKESPAN optimization problem consists in defining  $\mu$  (task mapping),  $\text{st}$  (task starting times) and  $\text{com}$  (communication starting times) so that the makespan  $M$  defined in Equation (1) is minimized.*

Note that this classical scheduling problem is NP-complete, even without communications, since the problem with  $n$  weighted independent tasks and  $p = 2$  processors is equivalent to the 2-partition problem [6].

### IV. ALGORITHMS

We propose novel heuristic approaches to solve the MINMAKESPAN problem, using a recent directed graph partitioner [11]. We compare the results with classical list-based

		(Cluster-based only)			task priority type	(partition priority) task priority	placement
		clustering approach	load balancing	produced clusters			
List-based	BL-EST				static	BL	EST-processor
	ETF				dynamic	EST	
Clustering-based	DSC-GLB-ETF	cyclic cluster graph limited refinement	Guided Load Balancing	Cyclic, non-convex graph	dynamic	TL+BL	EST task within cluster (processor)
<i>Proposed</i> Partitioning-based	*-PART	acyclic cluster graph better refinement		Directed Acyclic Graph	priority type of * (static or dynamic)	priority of * (BL or EST)	EST-Processor
	*-BUSY						EST-idle Processor
	*-MACRO						Earliest (Part) Finish Time-Processor

Table II: Heuristic approaches to solve MINMAKESPAN.

and clustering heuristics, that we first describe and adapt for the duplex single-port communication model (Section IV-A). Next, we introduce three variants of partition-assisted list-based scheduling heuristics in Section IV-B.

For convenience, Table I summarizes acronyms used in the paper, in particular in the heuristic names, and Table II summarizes the main features of all considered approaches.

#### A. State-of-the-art scheduling heuristics

We first consider the best alternatives from the list-based and cluster-based scheduling heuristics presented by Wang and Sinnen [27]. We consider one static list-scheduling heuristic (BL-EST), the best dynamic priority list-based scheduling heuristic for real application graphs (ETF), and the best cluster-based scheduling heuristic (DSC-GLB-ETF).

a) BL-EST: This simple heuristic maintains an ordered list of *ready* tasks, i.e., tasks that can be executed since all their predecessors have already been executed. Let  $\text{Ex}$  be the set of tasks that have already been executed, and let  $\text{Ready}$  be the set of ready tasks. Initially,  $\text{Ex} = \emptyset$ , and  $\text{Ready} = \{v_i \in V \mid \text{Pred}[v_i] = \emptyset\}$ . Once a task has been executed, new tasks may become ready. At any time, we have:

$$\text{Ready} = \{v_i \in V \setminus \text{Ex} \mid \text{Pred}[v_i] = \emptyset \text{ or } \forall (v_j, v_i) \in E, v_j \in \text{Ex}\}. \quad (2)$$

In the first phase, tasks are assigned a *priority*, which is designated to be its bottom level (hence the name BL). The bottom level  $\text{bl}(i)$  of a task  $v_i \in V$  is defined as the largest weight of a path from  $v_i$  to a target node (vertex without successors), including the weight  $w_i$  of  $v_i$ , and all communication costs. Formally,

$$\text{bl}(i) = w_i + \begin{cases} 0 & \text{if } \text{Succ}[v_i] = \emptyset; \\ \max_{v_j \in \text{Succ}[v_i]} c_{i,j} + \text{bl}(j) & \text{otherwise.} \end{cases} \quad (3)$$

In the second phase, tasks are assigned to processors. At each iteration, the task of the  $\text{Ready}$  set with the highest priority is selected and scheduled on the processor that would result in the earliest start time of that task. The start time depends on the time when that processor becomes available, the communication costs of its input edges, and the finish time of its predecessors. We keep track of the finish time of each processor  $P_k$  ( $\text{comp}_k$ ), as well as the finish time of sending ( $\text{send}_k$ ) and receiving ( $\text{recv}_k$ ) operations. When we tentatively schedule a task on a processor, if several communications are needed (meaning that at least two predecessors of the task are mapped on other processors), they cannot be performed at the same time with the duplex single-port communication model. The communications from the predecessors

are, then, performed as soon as possible (respecting the finish time of the predecessor and the available time of the sending and receiving ports) in the order of the finish time of the predecessors.

This heuristic is called BL-EST, for *Bottom-Level Earliest-Start-Time*, and is described in Algorithm 1. The  $\text{Ready}$  set is stored in a max-heap structure for efficiently retrieving the tasks with the highest priority, and it is initialized at line 3. The computation of the bottom levels for all tasks (line 1) can easily be performed in a single traversal of the graph in  $O(|V| + |E|)$  time, see for instance [17]. The main loop traverses the DAG and tentatively schedules a task with the largest bottom level on each processor in the loop lines 7-13. The processor with the earliest start time is then saved, and all variables are updated on lines 14-23. When updating  $\text{com}(j, i)$ , if  $v_i$  and its predecessor  $v_j$  are mapped on the same processor, communication start time is artificially set to  $\text{st}(j) + w_j - c_{j,i}$  in line 17, so that  $\text{st}(i)$  can be computed correctly in line 21. Finally, the list of ready tasks is updated line 23, i.e.,  $\text{Ex} \leftarrow \text{Ex} \cup \{v_i\}$ , and new ready tasks according to Equation (2) are inserted into the max-heap.

The total time complexity of Algorithm 1 is hence  $O(p^2|V| + |V| \log |V| + p|E|)$ :  $p^2|V|$  for lines 10-13 (for each processor, we need to keep and update temporary send/receive arrays),  $|V| \log |V|$  for the heap operations (we perform  $|V|$  times the extraction of the maximum, and the insertion of new ready tasks into the heap), and  $p|E|$  for lines 11-12. The space complexity is  $O(p + |V| + |E|)$ .

b) ETF: We also consider a dynamic priority list scheduler, ETF. For each ready task, this algorithm computes the earliest start time (EST) of the task. Then, it schedules the ready task with the earliest EST, hence the name ETF, for *Earliest EST First*. Since we tentatively schedule each ready task, the time complexity of ETF is higher than BL-EST; it becomes  $O(p^2|V|^2 + p|V||E|)$ . The space complexity is the same as BL-EST, i.e.,  $O(p + |V| + |E|)$ .

c) DSC-GLB-ETF: The clustering scheduling algorithm used as a basis for comparison is one of the best ones identified by Wang and Sinnen [27], namely, the DSC-GLB-ETF algorithm. It uses dominant sequence clustering (DSC), then merges clusters with guided load balancing (GLB), and finally orders tasks using earliest EST first (ETF). We refer the reader to [27] for more details about this algorithm.

#### B. Partition-based heuristics

The partition-based heuristics start by computing an acyclic partition of the DAG, using a recent DAG partitioner [11]. This



---

**Algorithm 1:** BL-EST algorithm

---

**Data:** Directed graph  $G = (V, E)$ , number of processors  $p$   
**Result:** For each task  $v_i \in V$ , allocation  $\mu(i)$  and start time  $\text{st}(i)$ ; For each  $(v_i, v_j) \in E$ , start time  $\text{com}(i, j)$

```
1  $\text{bl} \leftarrow \text{ComputeBottomLevels}(G)$ 
2  $\text{Ready} \leftarrow \text{EmptyHeap}$ 
3 Insert  $v_i$  in Ready with key  $\text{bl}(i)$  for all  $v_i$  without any
  predecessors
4 for  $k = 1$  to  $p$  do
5    $\text{comp}_k \leftarrow 0$ ;  $\text{send}_k \leftarrow 0$ ;  $\text{recv}_k \leftarrow 0$ ;
6 while Ready is not empty do
7    $v_i \leftarrow \text{extractMax}(\text{Ready})$ 
8   Sort  $\text{Pred}[v_i]$  in a non-decreasing order of the finish times
9   for  $k = 1$  to  $p$  do
10     $\text{begin}_k \leftarrow \text{comp}_k$ 
11    for  $v_j \in \text{Pred}[v_i]$  do
12      Update earliest possible begin time  $\text{begin}_k$  with
        the latest finishing predecessor communication.
13   $k^* \leftarrow \text{argmin}_k \{\text{begin}_k\}$  // Best Processor
14   $\mu(i) \leftarrow k^*$ 
15   $\text{st}(i) \leftarrow \text{comp}_{k^*}$ 
16  for  $v_j \in \text{Pred}[v_i]$  do
17    if  $\mu(j) = k^*$  then  $\text{com}(j, i) \leftarrow \text{st}(j) + w_j - c_{j,i}$ 
18    else
19       $\text{com}(j, i) \leftarrow \max\{\text{st}(j) + w_j, \text{send}_{\mu(j)}, \text{recv}_{k^*}\}$ 
20       $\text{send}_{\mu(j)} \leftarrow \text{recv}_{k^*} \leftarrow \text{com}(j, i) + c_{j,i}$ 
21     $\text{st}(i) \leftarrow \max\{\text{st}(i), \text{com}(j, i) + c_{j,i}\}$ 
22   $\text{comp}_{k^*} \leftarrow \text{st}(i) + w_i$ 
23 Insert new ready tasks into Ready
```

---

acyclic DAG partitioner takes a DAG with vertex and edge weights, a number of parts  $K$ , and an allowable imbalance parameter  $\varepsilon$  as input. Its output is a partition of the vertices of  $G$  into  $K$  nonempty pairwise disjoint and collectively exhaustive parts satisfying three conditions: (i) the weight of the parts are balanced, i.e., each part has a total vertex weight of at most  $(1 + \varepsilon) \frac{\sum_{v_i \in V} w_i}{K}$ ; (ii) the edge cut is minimized; (iii) the partition is acyclic; in other words, the inter-part edges between the vertices from different parts should preserve an acyclic dependency structure among the parts. We use this tool to partition the task graph into  $K = \alpha \times p$  parts, where  $\alpha \geq 1$  can be interpreted as the average number of clusters per processor. We choose an imbalance parameter of  $\varepsilon = 1.1$  to have relatively balanced clusters; other values of  $\varepsilon$  led to similar results. It may not always be possible to find a feasible partition with the given constraints, especially for small graphs and large  $\alpha$  and  $K$  values. However, since our main goal is to achieve good clustering, not perfect balance, we will continue with whatever partitioning found by our tool, even if it is not balanced (which only happened very rarely in our experiments).

Given  $K$  parts  $V_1, \dots, V_K$  forming a partition of the DAG, we propose three variants of scheduling heuristics. Note that the variants are designed on top of BL-EST and ETF, but they can easily be adapted to any other list-based scheduling algorithm since, in essence, these heuristics are capturing a hy-

brid approach between cluster-based and list-based scheduling algorithms using DAG partitioning.

a) \*-PART: The first variant, denoted \*-PART, is used in this paper on top of BL-EST or ETF. The BL-EST-PART heuristic (resp. ETF-PART) performs a list scheduling heuristic similar to BL-EST described in Algorithm 1 (resp. similar to ETF), but with the additional constraint that two tasks that belong to the same part must be mapped on the same processor. This means that once a task of a part has been mapped, we enforce that other tasks of the same part share the same processor, and hence do not incur any communication cost among the tasks of the same part. The pseudo-code of \*-PART can be found in the companion research report [19]. The time complexity of the partitioner is linear on the graph size [11], hence the complexity of BL-EST-PART (resp. ETF-PART) is the same as BL-EST (resp. ETF).

b) \*-BUSY: One drawback of the \*-PART heuristics is that it may happen that the next ready task is in a part that we are just starting (say  $V_\ell$ ), while some other parts have not been entirely scheduled. For instance, if processor  $P_j$  has already started processing a part  $V_{\ell'}$  but has not scheduled all of the tasks of  $V_{\ell'}$  yet, \*-PART may decide to schedule the new task from  $V_\ell$  onto the same processor if it will start at the earliest time. This may overload the processor and delay other tasks from both  $V_{\ell'}$  and  $V_\ell$ .

The second variant, \*-BUSY, checks whether a processor is already busy with an on-going part, and it does not allocate a ready task from another part to a busy processor, unless if all processors are busy. In this latter case, \*-BUSY behaves similarly to \*-PART. The pseudo-code of \*-BUSY can be found in the companion research report [19], and the complexity of this variant is the same as the list scheduling heuristic on top of which the variant is run, in our case BL-EST or ETF.

c) \*-MACRO: The last variant, \*-MACRO, further focuses on the parts, and schedules a whole part before moving to the next one, so as to avoid problems discussed earlier. This heuristic strongly relies on the fact that the partitioning is acyclic, and hence it is possible to process parts one after another in a topological order.

We extend the definition of ready tasks to parts. A part is ready if all its predecessor parts have already been processed. Hence, when a part is ready, all predecessors of tasks from that part have already been scheduled. We also extend the definition of bottom level to parts, by taking the maximum bottom level of tasks in the part.

The generic \*-MACRO is detailed in Algorithm 2. The algorithm relies on two priority algorithms, one for selecting parts, and one for selecting tasks. These priorities can be static, such as BL (selects parts or tasks with maximum bottom level), or dynamic, such as earliest start time as in ETF. Once a part has been selected, the algorithm tentatively schedules the whole part on each processor (lines 4-14). Tasks within the part are selected with the second priority algorithm. Incoming communications are scheduled at that time to ensure the

**Algorithm 2:** \*-MACRO algorithm

---

**Data:** Directed graph  $G = (V, E)$ , number of processors  $p$ , acyclic partition of  $G$ :  $V_1, \dots, V_K$ , a partition priority algorithm  $PP$ , a task priority algorithm  $TP$

**Result:** For each task  $v_i \in V$ , allocation  $\mu(i)$  and start time  $st(i)$ ; For each  $(v_i, v_j) \in E$ , start time  $com(i, j)$

```

1 Initialize ReadyParts with all  $V_i$  without any predecessors
2 while ReadyParts is not empty do
3    $V_i \leftarrow PP(\text{ReadyParts})$ 
      //  $PP()$  returns highest priority part
      from the ReadyParts list.
4   for  $k = 1$  to  $p$  do
5      $end_k \leftarrow comp_k$ 
6     Initialize Ready with all tasks from  $V_i$  with no
       unscheduled predecessors
7     while Ready is not empty do
8        $v_x \leftarrow TP(\text{Ready})$ 
          //  $TP()$  returns highest priority
          task from the Ready list.
9       Sort  $Pred[v_x]$  in a non-decreasing order of the
       finish times
10      Assign communication times (in  $Pred[v_x]$  order)
       and update computation times
11       $\mu(i) \leftarrow k$ 
12      Update  $st(x)$ ,  $com$  and  $comp$ 
13      Update  $end_k$  with the latest finishing task
14      Insert new ready tasks from same part into Ready
15    $k^* \leftarrow \text{argmin}_k \{end_k\}$  // Best Processor
16   Schedule part  $V_i$  to processor  $k^*$  (with the same procedure
    as in lines 6-14)
17   Insert new ready parts into ReadyParts

```

---

single-port model, and outgoing communications are left for later. The processor that minimizes the *finish time* is selected, and the part is assigned to this processor, since we aim at finishing a part as soon as possible to minimize the makespan. The finish times for computation, sending, and receiving are updated. Once a part has been scheduled entirely, the list of ready parts is updated, and the next ready part with highest priority is selected.

An instantiation of this algorithm with the bottom level priority algorithms (BL-MACRO) is available in the companion research report [19]. In ETF-MACRO, similarly to heuristic ETF, we tentatively schedule each ready part, and then each task, and at each step, we keep the best schedule. The time complexity of these variants are slightly different than the list scheduling heuristics on top of which the variant is run, because of part-by-part scheduling. For ETF-MACRO, the complexity is  $O(p^4 + p^3|V| + p^2|V|^2 + p|V||E|)$ .

## V. SIMULATION RESULTS

We first describe the simulation setup in Section V-A, in particular, the different instances that we use in the simulations. Next, we compare the baseline algorithms under different communication models (communication-delay model vs. realistic model) in Section V-B. Section V-C shows the impact of the number of parts used by the partitioner, the communication-to-computation ratio (CCR), the number of

Graph	$ V $	$ E $	Degree		#source	#target
			max.	avg.		
598a	110,971	741,934	26	13.38	6,485	8,344
caidaRouterLev.	192,244	609,066	1,071	6.34	7,791	87,577
delaunay-n17	131,072	393,176	17	6.00	17,111	10,082
email-EuAll	265,214	305,539	7,630	2.30	260,513	56,419
fe-ocean	143,437	409,593	6	5.78	40	861
ford2	100,196	222,246	29	4.44	6,276	7,822
luxembourg-osm	114,599	119,666	6	4.16	3,721	9,171
rgg-n-2-17-s0	131,072	728,753	28	5.56	598	615
usroads	129,164	165,435	7	2.56	6,173	6,040
vsp-mod2-pgp2.	101,364	389,368	1,901	7.68	21,748	44,896

Table III: Instances from the UFL Collection [3].

Graph	$ V $	$ E $	Degree		#source	#target
			max.	avg.		
cholesky	1,030,204	1,206,952	5,051	2.34	333,302	505,003
fibonacci	1,258,198	1,865,158	206	3.96	2	296,742
quicksort	1,970,281	2,758,390	5	2.80	197,030	3
RSBench	766,520	1,502,976	3,074	3.96	4	5
Smith-water.	58,406	83,842	7	2.88	164	6,885
UTS	781,831	2,061,099	9,727	5.28	2	25
XSbench	898,843	1,760,829	6,801	3.92	5	5

Table IV: Instances from OCR [30].

processors, and the edge cut. Finally, we present detailed simulation results in Section V-D and summarize these results in Section V-E.

The code is available at <http://tda.gatech.edu/software/dagPscheduler/> so that interested readers can instantiate their graphs and repeat simulations for reproducibility purpose.

## A. Simulation setup

The experiments were conducted on a computer equipped with dual 2.1 GHz Xeon E5-2683 processors and 512GB memory. We have performed an extensive evaluation of the proposed cluster-based scheduling heuristics on instances coming from three sources.

The first set of instances is from Wang and Sinnen's work [27]. This set contains roughly 1600 instances of graphs, each having 50 to 1151 nodes. All graphs have three versions for CCRs 0.1, 1, and 10. The dataset includes a wide range of real world, regular structure, and random structure graphs; more details about them are available in the original paper [27]. Since the graphs are up to 1151 nodes, we refer to this dataset as the *small* dataset.

The second set of instances is obtained from the matrices available in the SuiteSparse Matrix Collection (formerly known as the University of Florida Sparse Matrix Collection) [3]. From this collection, we picked ten matrices satisfying the following properties: listed as binary, square, and has at least 100000 rows and at most  $2^{26}$  nonzeros. For each such matrix, we took the strict upper triangular part as the associated DAG instance, whenever this part has more nonzeros than the lower triangular part; otherwise we took the lower triangular part. The ten graphs from the UFL dataset and their characteristics are listed in Table III.

The third set of instances is from the Open Community Runtime (OCR), an open source asynchronous many-task runtime that supports point-to-point synchronization and disjoint data blocks [30]. We use seven benchmarks from the OCR

repository<sup>1</sup>. These benchmarks are either scientific computing programs or mini-apps from real-world applications whose graphs' characteristics are listed in Table IV.

To cover a variety of applications, we consider UFL and OCR instances with random edge costs and random vertex weights, using different communication-to-computation ratios (CCRs). For a graph  $G = (V, E)$ , the CCR is formally defined as

$$\text{CCR} = \frac{\sum_{(v_i, v_j) \in E} c_{i,j}}{\sum_{v_i \in V} w_i}.$$

In order to create instances with a target CCR, we proceed in two steps: (i) we first randomly assign costs and weights between 1 and 10 to each edge and vertex, and then (ii) we scale the edge costs appropriately to yield the desired CCR.

Since the ETF algorithms have a complexity in  $O(p^2|V|^2 + p|V||E|)$ , they are not suited to million-node graphs that are included in the OCR and UFL datasets. Hence, we have selected a subset of OCR and UFL graphs, namely, graphs with 10k to 150k nodes, denoted as the *medium* dataset. The *big* dataset contains all graphs from Tables III and IV.

#### B. Communication-delay model vs. realistic model

Our goal is to compare the new heuristics with the best competitors from the literature [27]. We call them the baseline heuristics, as they represent the current state-of-the-art. We have access to executables of the original implementation [27]. However, these heuristics assume a pure communication-delay model, where communications can all happen at the same time, given that the task initiating the communications has finished its computation. Hence, there is no need to schedule the communications in this model.

In our work, we have assumed a more realistic, duplex single-port communication model. Thus, we cannot directly compare the new heuristics with the executables of the baseline heuristics. We have, therefore, implemented our own version of the baseline algorithms (BL-EST, ETF as best list-based and DSC-GLB-ETF as best cluster-based scheduler) with the communication delay model, and compared the resulting makespans with those of Wang and Sinnen's implementation, denoted as "ETF [W&S]", in an attempt to validate our implementations. We show the performance profiles in Figure 2 for this comparison. In the performance profiles, we plot the percentages of the instances in which a scheduling heuristic obtains a makespan on an instance that is no larger than  $\theta$  times the best makespan found by any heuristic for that instance [4]. Therefore, the higher a profile at a given  $\theta$ , the better a heuristic is. Results on Figure 2 confirm those presented by Wang and Sinnen: with low CCR (CCR=0.1 or CCR=1), DSC-GLB-ETF is worse than ETF (the higher the better). However, when the CCR increases, the performance of DSC-GLB-ETF also increases, and it surpasses ETF for CCR=10 at the end [27].

Figure 2 also shows that our implementation of ETF performs better than ETF [W&S]. This may be due to tie-breaking

in case of equal ordering condition, that we could not verify in detail since we had only the executables. Our implementation ETF is, thus, a fair competitor, since it turns out to be better than the existing implementation.

Next, we converted our implementation of these algorithms into duplex single-port model, as explained in Section IV, in order to establish the baseline to compare the proposed heuristics. Figure 3 shows the performance profiles of our three baseline heuristics on the *small* dataset. From these results, we see that DSC-GLB-ETF is not well suited for the realistic communication model, since it performs pretty badly in comparison to ETF. BL-EST is also slightly worse than ETF, but it has a lower time complexity.

#### C. Impact of number of parts, processors, CCR, and edge cut

Here, we evaluate the impact of number of parts in the partitioning phase, number of processors, CCR of datasets, and edge cut of the partitioner on the quality of the proposed heuristics.

Figure 4 depicts the relative performance of BL-EST-PART, BL-EST-BUSY, and BL-MACRO compared to BL-EST on the *big* dataset as a function of  $\alpha$  for different number of processors,  $p = \{2, 4, 8, 16, 32\}$ , and CCR=10. We set the number of parts  $K = \alpha \times p$  and we have  $\alpha = \{1, 2, 3, 4, 6, 8, 10, 12, 14, 16\}$ . As seen in the figure, except BL-MACRO on  $p = 32$  processors, the new algorithms perform better than the baseline BL-EST for all values of  $\alpha$  that we tested. Even for the worst case, that is, on 32 processors, BL-MACRO performs better or comparable to BL-EST, when  $\alpha \leq 4$ . Therefore, we recommend to select  $\alpha \leq 4$ .

As shown in the previous studies (e.g., [27]), performance of the scheduling algorithms vary significantly with different CCRs, and in particular, clustering-based schedulers perform better for high CCRs, i.e., when communications are more costly than computations. Figure 5 shows the performance of the heuristics on the *big* dataset with varying CCR, i.e., for CCR= $\{1, 5, 10, 20\}$ , and for  $p = \{2, 4, 8, 16, 32\}$ . The results are the average of all input instances using the best  $\alpha$  value, for  $\alpha = \{1, 2, 3, 4\}$ , for that instance.

As expected, similar to existing clustering-based schedulers, the proposed heuristics give significantly better results than the BL-EST baseline. For instance, when CCR=20, for all numbers of processors in the figure, all partitioning-based heuristics give at least 50% better makespans.

Comparing the relative performance of BL-EST-PART and BL-EST-BUSY across the sub-figures, one observes that BL-EST-PART and BL-EST-BUSY have more or less stable performance with the increasing number of processors. Note that the performance of BL-EST-PART and BL-EST-BUSY mostly depends on the value of CCR, but remains the same when the number of processors varies. BL-MACRO performs worse than the other two heuristics for small values of CCR with an increasing number of processors. However, for tested values of  $p$ , the performance of BL-MACRO improves as the CCR increases, and finally it outperforms all other heuristics on average when the CCR is large enough.

<sup>1</sup><https://xstack.exascale-tech.com/git/public/apps.git>



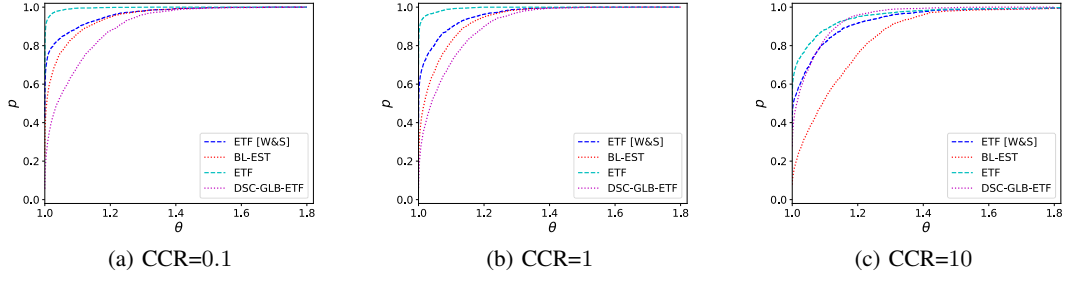


Figure 2: Performance profiles comparing our implementation of baseline heuristics with Wang and Sinnen's implementation of ETF, on the *small* data set, with the communication-delay model.

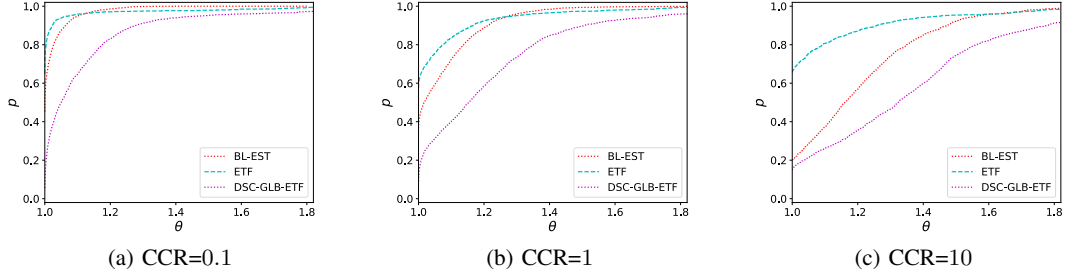


Figure 3: Performance profiles comparing baselines on the *small* dataset, with the duplex single-port communication model.

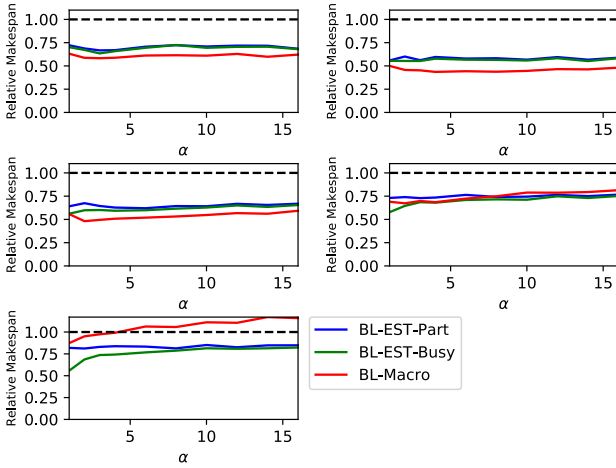


Figure 4: Relative makespan compared to BL-EST on the *big* dataset, as a function of the number of parts, with CCR=10 and with 2 (top left), 4 (top right), 8 (middle left), 16 (middle right), and 32 (bottom left) processors.

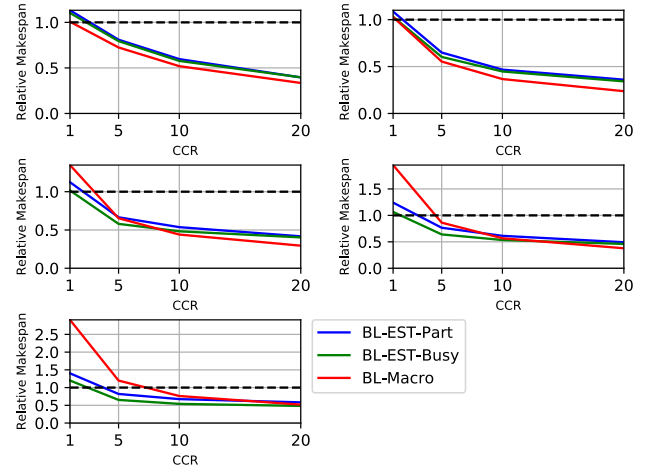


Figure 5: Relative makespan compared to BL-EST on the *big* dataset, as a function of the CCR, with 2 (top left), 4 (top right), 8 (middle left), 16 (middle right), and 32 (bottom left) processors.

We have carried out thorough experiments to see the effects of the edge cut of DAG partitioning in the final makespan. We observed that having a smaller edge cut in DAG partitioning yields a better makespan more than 82% of the time for all proposed heuristics, when the communication-to-computation ratio (CCR) is 10. Indeed, on the *small* dataset, we counted the instances where a better edge cut in partitioning gave a better makespan. Out of 9045 instances, there were 7494 such instances for \*-MACRO, 7519 for \*-PART, and 7477 for \*-BUSY, hence ranging between 82.6% and 83.1%.

#### D. Performance results and runtime comparison

We present the results on the *small*, *medium* and *big* datasets. We focus only on the BL-EST algorithm for the *big* dataset, since ETF does not scale well (due to quadratic time complexity on the number of vertices), and DSC-GLB-ETF shows poor results with the realistic communication model and smaller datasets. Let us consider xSBench graph as an example of how long it takes to run ETF on one of the *big* graphs. When we schedule this graph on two processors, the DAG partitioning algorithm runs in 9.5 seconds on average,

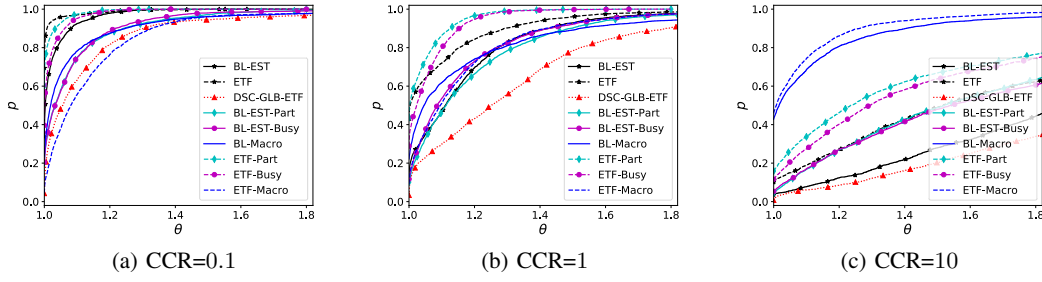


Figure 6: Performance profiles comparing all the algorithms on the *small* dataset with duplex single-port model.

Graph	CCR=1				CCR=20			
	BL-EST	BL-EST-PART	BL-EST-BUSY	BL-MACRO	BL-EST	BL-EST-PART	BL-EST-BUSY	BL-MACRO
598a	3058476	1.14	1.14	1.04	17038485	0.22	0.22	0.22
caidaRouterLevel	5337718	1.02	1.02	1.00	26745328	0.56	0.56	0.35
delaunay-n17	3606092	1.02	1.03	1.00	17567627	0.22	0.22	0.21
email-EuAll	7711619	1.00	1.00	0.98	67066585	0.19	0.19	0.19
fe-ocean	3949464	1.12	1.12	1.02	11573357	0.39	0.39	0.35
ford2	2781775	1.03	1.03	0.99	10538479	0.27	0.27	0.26
luxembourg-osm	3152973	1.01	1.01	1.00	4801062	0.66	0.66	0.66
rgg-n-2-17-s0	3601079	1.23	1.23	1.06	9094485	0.48	0.48	0.43
usroads	3550396	1.02	1.02	1.02	8428888	0.43	0.43	0.42
vsp-mod2-pgp2-slptsk	2794636	1.04	1.04	1.00	19887584	0.28	0.46	0.23
Cholesky	30603433	1.28	1.03	0.95	130153391	0.31	0.24	0.23
fibonacci	34601228	1.11	1.10	1.03	110167490	0.36	0.35	0.32
quicksort	54162227	1.01	1.01	1.00	173055640	0.32	0.32	0.31
RSBench	26941941	1.38	1.25	0.88	109245784	0.38	0.30	0.24
Smith-waterman	1661676	1.46	1.41	1.02	5694549	0.53	0.44	0.33
UTS	31904401	1.34	1.34	1.34	117598932	0.40	0.41	0.37
XSbench	41794985	1.15	1.15	1.02	77257208	0.64	0.63	0.60
<b>Geomean</b>	<b>1.00</b>	<b>1.13</b>	<b>1.11</b>	<b>1.02</b>	<b>1.00</b>	<b>0.36</b>	<b>0.36</b>	<b>0.32</b>

Table V: The makespan of BL-EST in absolute numbers, and those of BL-EST-PART, BL-EST-BUSY, and BL-MACRO relative to BL-EST on *big* dataset, when the number of processors  $p$  is 2, and for  $CCR=\{1, 20\}$ .

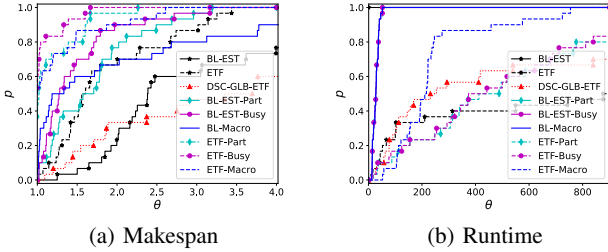


Figure 7: Performance profiles on *medium* dataset, with  $CCR=10$ .

and BL-EST-PART, BL-EST-BUSY, and BL-MACRO heuristics run under half a second, making the total time approximately 10 seconds. However, ETF algorithm takes 4759 seconds. On four processors, it goes up to 7507 seconds, and so on.

a) *Small dataset*: Figure 6 shows the comparison of all heuristics on the *small* dataset for  $CCR=\{0.1, 1, 10\}$ . While ETF remains the best with a small  $CCR=0.1$ , ETF-PART becomes better as soon as  $CCR=1$ . Finally, the performance of BL-MACRO and ETF-MACRO is striking for  $CCR=10$ , where the \*-MACRO variant clearly outperforms all other heuristics.

As seen before, DSC-GLB-ETF performs poorly in this case, since it is not designed for realistic duplex single-port communication model.

b) *Medium dataset*: Figure 7a shows the performance profile for the *medium* dataset. As expected, dynamic scheduling technique ETF and our ETF-based heuristics perform better than their BL-EST counterparts, as for the *small* dataset. Note that our heuristics perform better than the original versions they are built upon.

ETF and ETF-based algorithms' quality comes with a downside of high time complexity and consequently, slower algorithms due to their dynamic nature. Figure 7b shows runtime performance profiles. It is therefore the fraction of instances in which an algorithm gave a runtime no worse than the fastest algorithm, hence the higher the better. As expected, the static BL-EST approach runs much faster than dynamic approaches. DAG-partitioning introduces an overhead to proposed heuristics, but this is still negligible compared to the time complexity of the algorithms. BL-EST-PART, BL-EST-BUSY, and BL-MACRO heuristics also perform comparably fast even with partitioning time overhead. ETF and ETF-based heuristics run two to three orders of magnitude slower, making them infeasible to run on bigger graphs.

c) *Big dataset*: Table V displays the detailed results on the *big* dataset, with two processors, for  $CCR=1$  and  $CCR=20$ . Results for  $CCR=5$  and  $CCR=10$  are available in the companion research report [19]. On average, BL-EST-BUSY provides slightly better results than BL-EST-PART. When  $CCR=1$ , the heuristics often return a makespan that is

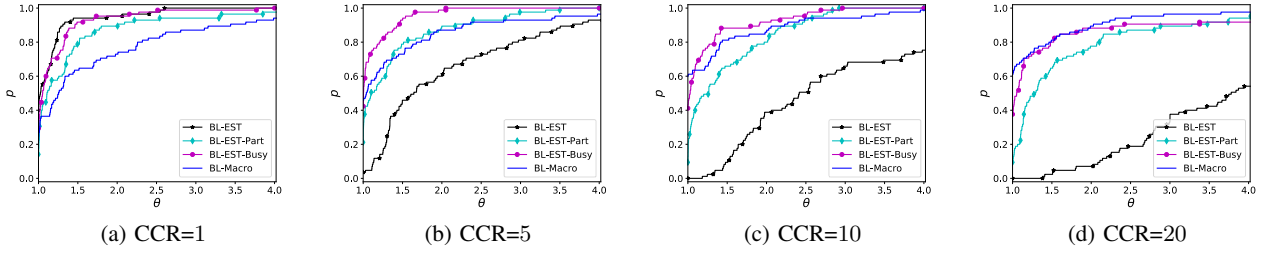


Figure 8: Performance profiles on *big* dataset, with  $CCR=\{1, 5, 10, 20\}$ .

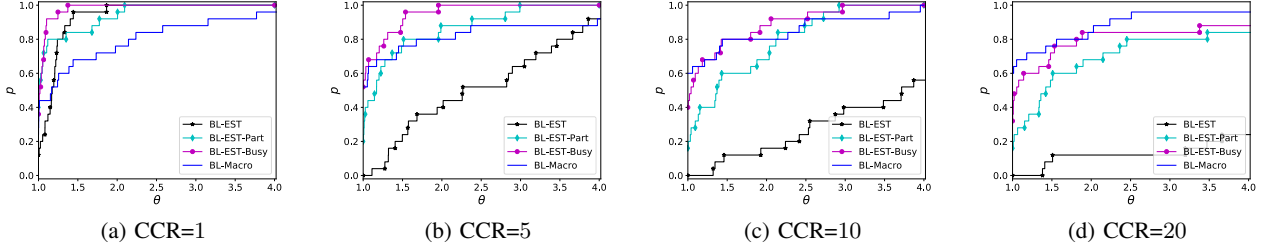


Figure 9: Performance profiles on *big* dataset when more than 10% of nodes are sources, with  $CCR=\{1, 5, 10, 20\}$ .

slightly larger than the one from BL-EST, on average by 13%, 11%, and 2%, respectively. However, for  $CCR=20$ , BL-EST-PART and BL-EST-BUSY obtain a makespan 2.8 times smaller than the baseline, and it goes up to 3.1 times smaller for BL-MACRO, when we average results on the whole *big* dataset and considering an architecture with two processors.

Figure 8 shows the performance profiles of these four algorithms for  $CCR=\{1, 5, 10, 20\}$ . When  $CCR=1$ , BL-EST performs best but BL-EST-BUSY performs very close to it. However, when the value of CCR is increasing, it is more and more important to handle communications correctly. We observe that the proposed three heuristics outperform the baseline (BL-EST) as the CCR increases. When  $CCR=5$ , in about 90% of all cases, BL-EST-BUSY’s makespan is within  $1.5\times$  of the best result, whereas this fraction is only 40% for BL-EST. Starting with  $CCR=10$ , the proposed heuristics completely dominate BL-EST algorithm. For all values of CCRs, BL-EST-BUSY outperforms BL-EST-PART. However, BL-MACRO performs worse than BL-EST-PART and BL-EST-BUSY when  $CCR=1$ , and gradually outperforms the other two as the CCR increases.

To understand the nature of datasets where the proposed heuristics and the baseline behave differently, we divided the *big* dataset into two subsets, the graphs consisting of more than 10% of the nodes as sources, and the ones with less than 10%. Figure 9 shows how the algorithms’ quality differ for when the DAGs have more than 10% nodes as sources. With a lot of sources, BL-EST baseline performs badly while BL-MACRO performs better compared to the case with fewer sources. This is due to the inherent nature of DAG-partitioning followed by cluster-by-cluster scheduling. Consider a DAG of clusters with one source cluster. BL-MACRO would need to schedule all of the nodes in this cluster in one processor to

start utilizing any other processor available. When the number of source clusters is high, this heuristic can start efficiently using more processors right from the start.

#### E. Summary

Overall, the proposed meta-heuristics significantly improve the makespan found by the baseline heuristics they are applied on, as empirically shown with a wide range of graph instances. The results confirm the correlation between the edge cut found during the partitioning phase and the makespan at the end. The benefit of a good partitioning with minimum edge cut objective shows itself clearly, especially when the CCR is high.

The results show that \*-PART and \*-BUSY behave consistently, and that they provide a steady improvement over the baselines. Furthermore, their relative performance (compared to the baseline) does not depend on the number of processors, which means that these heuristics scale well. They perform even better when the ratio between communication and computation is large.

The \*-MACRO’s performance has a higher variance. This meta-heuristic tries to have more of a “global” view during scheduling, by tentatively scheduling whole parts instead of deciding the mapping when it is only at the first node of the part and dictating the rest (as done by \*-PART and \*-BUSY). It seems to not scale when the number of processors increases. Nevertheless, when the ratio between communication and computation is large, it usually outperforms all the other heuristics. Also, the experiments show that when the input instance to be scheduled has higher percentage of sources (source parts), \*-MACRO is even more likely to outperform other heuristics.

## VI. CONCLUSION

We proposed three new partitioning-assisted list-based scheduling techniques (or meta-heuristics) based on an acyclic partition of the DAGs: \*-PART, \*-BUSY, and \*-MACRO. The acyclicity of the partition ensures that we can schedule a part in its entirety as soon as its input nodes are available. Hence, we have been able to design specific list-scheduling techniques that would not have been possible without an acyclic partition. To the best of our knowledge, this is the first partitioning-assisted list-scheduler using a multi-level *directed* DAG partitioner for the clustering phase. The acyclicity is well suited to identify data locality in the DAG, and it allows the design of specific allocation strategies, such as \*-MACRO. The proposed meta-heuristics are generic and can be combined with any classical list-scheduling heuristic, and used with any acyclic partitioner.

We compared our scheduling techniques with the widely used BL-EST, ETF, and DSC-GLB-ETF heuristics, adapted to the realistic duplex single-port communication model. The results are striking, with the new heuristics consistently improving the makespan. Even though \*-MACRO does not seem to scale well with the number of processors, it delivers the best results in several cases, while \*-PART and \*-BUSY are consistently good. For instance, the proposed \*-PART (resp. \*-BUSY and \*-MACRO) algorithms achieve a makespan 2.6 (resp. 3.1 and 3.3) times smaller than BL-EST when considering the *big* dataset with  $CCR = 20$ , averaging over all processor numbers. Furthermore, if we pick the best of the three heuristics for each instance, it is four times better.

As future work, we plan to consider *convex* partitioning instead of acyclic partitioning, which is less restrictive and hence exposes more parallelism. To the best of our knowledge, there is no top-down convex partitioning technique available, which we plan to investigate. Also, an adaptation of the proposed heuristics to heterogeneous processing systems can be carried out.

*Acknowledgment:* We thank Oliver Sinnen for providing us the Java binaries of their implementation and the datasets they used in their studies [27], and the referees for valuable feedback.

## REFERENCES

- [1] T. L. Adam, K. M. Chandy, and J. Dickson, "A comparison of list schedules for parallel processing systems," *Communications of the ACM*, vol. 17, no. 12, pp. 685–690, 1974.
- [2] I. Ahmad and Y.-K. Kwok, "On exploiting task duplication in parallel program scheduling," *IEEE T. Parall. Distr.*, vol. 9, no. 9, pp. 872–892, 1998.
- [3] T. A. Davis and Y. Hu, "The University of Florida sparse matrix collection," *ACM Trans. Math. Softw.*, vol. 38, no. 1, pp. 1:1–1:25, 2011.
- [4] E. D. Dolan and J. J. Moré, "Benchmarking optimization software with performance profiles," *Mathematical programming*, vol. 91, no. 2, pp. 201–213, 2002.
- [5] I. T. Foster, M. Fidler, A. Roy, V. Sander, and L. Winkler, "End-to-end quality of service for high-end applications," *Computer Communications*, vol. 27, no. 14, pp. 1375–1388, 2004.
- [6] M. R. Garey and D. S. Johnson, *Computers and Intractability, a Guide to the Theory of NP-Completeness*. W.H. Freeman and Company, 1979.
- [7] A. Gerasoulis and T. Yang, "On the granularity and clustering of directed acyclic task graphs," *IEEE T. Parall. Distr.*, vol. 4, no. 6, pp. 686–701, June 1993.
- [8] T. Glatard, J. Montagnat, D. Lingrand, and X. Pennec, "Flexible and efficient workflow deployment of data-intensive applications on grids with MOTEUR," *Int. Journal of High Performance Computing and Applications*, 2008.
- [9] R. L. Graham, "Bounds on multiprocessing timing anomalies," *SIAM Journal on Applied Mathematics*, vol. 17, no. 2, pp. 416–429, 1969.
- [10] C. Hanen and A. Munier, "An approximation algorithm for scheduling dependent tasks on  $m$  processors with small communication delays," in *Proceedings 1995 INRIA/IEEE Symposium on Emerging Technologies and Factory Automation. ETFA'95*, vol. 1, Oct 1995, pp. 167–189 vol.1.
- [11] J. Herrmann, M. Y. Özkaya, B. Uçar, K. Kaya, and Ü. V. Çatalyürek, "Acyclic partitioning of large directed acyclic graphs," Inria - Research Centre Grenoble – Rhône-Alpes, Research Report RR-9163, Mar 2018. [Online]. Available: <https://hal.inria.fr/hal-01744603>
- [12] J.-J. Hwang, Y.-C. Chow, F. D. Anger, and C.-Y. Lee, "Scheduling precedence graphs in systems with interprocessor communication times," *SIAM Journal on Computing*, vol. 18, no. 2, pp. 244–257, 1989.
- [13] K. Jansen, F. Land, and M. Kaluza, "Precedence scheduling with unit execution time is equivalent to parametrized biclique," in *SOFSEM 2016: Theory and Practice of Computer Science*, R. M. Freivalds, G. Engels, and B. Catania, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2016, pp. 329–343.
- [14] H. Kanemitsu, M. Hanada, and H. Nakazato, "Clustering-based task scheduling in a large number of heterogeneous processors," *IEEE T. Parall. Distr.*, vol. 27, no. 11, pp. 3144–3157, Nov 2016.
- [15] Y.-K. Kwok and I. Ahmad, "Dynamic critical-path scheduling: An effective technique for allocating task graphs to multiprocessors," *IEEE T. Parall. Distr.*, vol. 7, no. 5, pp. 506–521, 1996.
- [16] —, "Benchmarking and comparison of the task graph scheduling algorithms," *Journal of Parallel and Distributed Computing*, vol. 59, no. 3, pp. 381–422, 1999.
- [17] —, "Static scheduling algorithms for allocating directed task graphs to multiprocessors," *ACM Comput. Surv.*, vol. 31, no. 4, pp. 406–471, 1999.
- [18] S. Mingsheng, S. Shixin, and W. Qingxian, "An efficient parallel scheduling algorithm of dependent task graphs," in *Proc. of 4th Int. Conf. on Parallel and Distributed Computing, Applications and Technologies, PDCAT*. IEEE, 2003, pp. 595–598.
- [19] M. Y. Özkaya, A. Benoit, B. Uçar, J. Herrmann, and Ü. V. Çatalyürek, "A scalable clustering-based task scheduler for homogeneous processors using DAG partitioning," Inria - Research Centre Grenoble – Rhône-Alpes, Tech. Rep. RR-9185, 2018, available online at [hal.inria.fr](https://hal.inria.fr).
- [20] C. Picouleau, "New complexity results on scheduling with small communication delays," *Discrete Applied Mathematics*, vol. 60, no. 1, pp. 331 – 342, 1995.
- [21] A. Radulescu and A. J. Van Gemund, "Low-cost task scheduling for distributed-memory machines," *IEEE T. Parall. Distr.*, vol. 13, no. 6, pp. 648–658, 2002.
- [22] V. Sarkar, "Partitioning and scheduling parallel programs for execution on multiprocessors," Stanford Univ., CA (USA), Tech. Rep., 1987.
- [23] G. C. Sih and E. A. Lee, "A compile-time scheduling heuristic for interconnection-constrained heterogeneous processor architectures," *IEEE T. Parall. Distr.*, vol. 4, no. 2, pp. 175–187, 1993.
- [24] O. Sinnen, *Task Scheduling for Parallel Systems*. New York, NY, USA: Wiley Series on Par. and Distr. Computing, Wiley-Interscience, 2007.
- [25] O. Svensson, "Conditional hardness of precedence constrained scheduling on identical machines," in *Proceedings of the Forty-second ACM Symposium on Theory of Computing*, ser. STOC'10. New York, NY, USA: ACM, 2010, pp. 745–754.
- [26] H. Topcuoglu, S. Hariri, and M. Y. Wu, "Performance-effective and low-complexity task scheduling for heterogeneous computing," *IEEE T. Parall. Distr.*, vol. 13, no. 3, pp. 260–274, 2002.
- [27] H. Wang and O. Sinnen, "List-scheduling vs. cluster-scheduling," *IEEE T. Parall. Distr.*, 2018, in press.
- [28] M.-Y. Wu and D. D. Gajski, "Hypertool: A programming aid for message-passing systems," *IEEE T. Parall. Distr.*, vol. 1, no. 3, pp. 330–343, 1990.
- [29] T. Yang and A. Gerasoulis, "DSC: Scheduling parallel tasks on an unbounded number of processors," *IEEE T. Parall. Distr.*, vol. 5, no. 9, pp. 951–967, 1994.
- [30] L. Yu and V. Sarkar, "GT-Race: Graph traversal based data race detection for asynchronous many-task runtimes," in *Euro-Par 2018: Parallel Processing*. Springer, 2018.