

# Comparison of Ranking Methods in HEFT Scheduling

Ryan Bunney  
email: 21298244@student.uwa.edu.au

May 2017

## Abstract

Efficient scheduling is critical to the performance of many problems faced in day-to-day life. The task scheduling problem has been demonstrated to be NP-Complete in both general and restricted cases. This has resulted in a large body of literature that details the variety of efficient approaches to creating good quality schedules. The context of many task-scheduling problems exist within heterogeneous computing environments; the Heterogeneous Earliest-Finish-Time (HEFT) algorithm is a popular list-scheduling algorithm that has achieved positive response for its good quality schedules and relatively-efficient execution times in heterogeneous system. Central to the HEFT algorithm is the use of the upward rank heuristic, which ranks each task based on its respective communication and computation costs. This method generates a sorted list that is a unique topological order of tasks. The original HEFT paper provides no validation for why this ranking method generates a better sort than a standard topological sort, which makes it difficult for researchers to determine the impact the rank has on the final schedule, and what areas of the HEFT algorithm can be improved or extended to increase its performance. An implementation of HEFT that sorts tasks using both the upward rank and a topological sort was written to provide a comparative analysis of the effect of each sorting method on the final schedule length and execution time of the algorithm. The comparison study, based on randomly generated graphs of increasing size, shows that the ranking sort is a significant improvement over a standard topological sort, in both the schedule length and execution time of the overall algorithm.

**Keywords:** Scheduling, Workflow, Heterogeneous, Distributed computing.

# 1 Introduction

The scheduling problem is a well-established and extensively-researched topic within computer science [5, 6, 10, 14]. Applications of scheduling include project workflow scheduling, air-traffic control, and big data processing – to name a few [19]. Scheduling has been found to be an NP-Complete problem, in the general sense; as a result, most research is dedicated to developing different heuristics to determine schedules for particular problems. In the case of multi-processor scheduling problems, in which the problem is focused on scheduling computational tasks onto many processors, the aim of the algorithm is to reduce the time it takes for all the tasks to be completed, and utilise the resources available as effectively as possible [8]. Some approaches work on the assumption that system resources are homogeneous [5, 9]; others attempt to utilise the extra resources that might be available on heterogeneous processors [10, 16]. The Heterogeneous Earliest-Finish-Time (HEFT) algorithm is a popular scheduling heuristic that has achieved a lot of attention for its performance and schedule quality [1]. HEFT is split into two parts; task-prioritisation and processor selection. The first part, task prioritisation is of particular interest to those interested in improving the quality of the resultant schedule. The prioritisation method used in HEFT sorts each task in order of priority generated by a ranking heuristic, before they are then assigned to a processor. This ranking method retains the precedence requirements of each task, and can be demonstrated to be a topological order of the task graph [16].

The original HEFT paper provides no justification for why its ranking method provides a better order of tasks for computation than standard a topological sort would. We believe this is an important oversight; if one is interested in improving the effectiveness of HEFT for various applications[1, 11], it is necessary to understand why certain methods perform better than others within the algorithm. If the ranking heuristic is not better than a generic topological sort, then time spent on perfecting the algorithm with respect to the heuristic could be utilised elsewhere. We aim to demonstrate that the ranking heuristic is a better method, and posit that when compared to a non-unique topological sort, the HEFT ranking algorithm produces a better sort that results in shorter schedule lengths.

This paper tests the effectiveness of the HEFT upward-rank heuristic against a topological sort of the task graph. We base our assessment on two criteria; firstly, the length of the resultant schedule, referred to as the *makespan* in the literature [17]; and secondly, the execution time of the algorithm. The experiments are conducted over a range of randomly generated graphs of different sizes. The comparison shows that both the makespan and

the execution time length of the HEFT algorithm is significantly lower when using the Rank heuristic demonstrated in [16]. We also notice a steeper increase in both execution time and schedule length as the number of nodes increases for a topological sort, and discuss potential reasons for this divergence.

The remainder of this paper is as follows: In the next section, the literature surrounding Task scheduling is summarised, and motivations for testing HEFT presented. In Section 3, the scheduling problem is formally defined, and an outline of the HEFT algorithm, ranking method, and processors selection presented. Section 4 provides detail on the implementation of the algorithm and the testing environment, and Section 5 presents a comparative study between the ranking and topological sorting approaches. In Section 6, we summarise our results and discuss the implications for future research.

## 2 Previous and related work

Scheduling problems are, in the general sense, NP-complete [3, 6, 10, 14, 17]. Kwok & Ahmad [10] note that only three sub-problems have polynomial time algorithms. As a result, algorithms and systems that aim to provide schedules for tasks must experiment with different methods that will find the closest-to-optimal makespan for a given set of tasks.

The scheduling problem is represented in the literature as a Directed Acyclic Graph (DAG) [5, 6, 17, 10, 6]. A sequence of tasks may have precedence constraints in which some tasks cannot start before others have completed, and in schedules with a definitive start and ending, tasks do not form loops [9]. The scheduling problem is formally defined in Section 3. In the follow sections, we summarise the areas of research that relate to the HEFT algorithm, and where it sits in the body of literature.

### 2.1 Homogeneous vs. Heterogeneous Algorithms

As identified in Kwok & Ahmad [10], research on scheduling has focused on scheduling tasks in homogeneous systems. We define a homogeneous system as a system in which computational capacity of resources in the system is equal across the board; this could include homogeneous processors on a CPU; or a HPC grid in which allocated resources are identical in capacity [10]. Early work on scheduling included focus on operating system and process scheduling for hardware [12], which are simplified to homogeneous systems. Seminal approaches to scheduling include the partitioning and scheduling model outlined by Sarkar [14], in which graphs are initially partitioned into

effective node ‘clusters’ and then scheduled based on the partitions. As identified in [5], this load-balancing technique favours homogeneous systems, as the complexity of partitioning a homogeneous graph is already a challenging problem [17].

As a result of the complexity of heterogeneous scheduling problems [6], less focus has been dedicated to heterogeneous research [10], and methods that are present in the literature are confronted by high scheduling costs [16]. HEFT has emerged as one of the most used techniques for scheduling on heterogeneous systems. Comparison studies against over 20 other algorithms have resulted in HEFT being described as one of the best scheduling techniques for schedule length [4] on heterogeneous systems.

## 2.2 Scheduling Techniques

The literature presents four common techniques for approaching the DAG scheduling problem; list scheduling, clustering, task duplication, and guided random search [16]. Of these, list scheduling is the most common due to its simple approach of allocating a list tasks to processors [10] and their generally superior performance [18]. Clustering algorithms are analogous to load balancing partitions in scheduling, and are typically intended for systems with an unbounded number of resources; this makes it less useful when applied to heterogeneous systems [5, 12]. Task duplication methods have been demonstrated as impractical for both homogeneous and heterogeneous contexts, due to their high complexity [10]. The most popular approach to random search methods is the use of Genetic Algorithms (GA). These have been shown to develop schedules of equal or better quality to that of a list scheduling heuristic [7]; however, they are let down by high execution times that outweigh the increase in schedule quality [16].

HEFT is a list-scheduling algorithm; it sorts the tasks into a list based on the ranking heuristic mentioned in Section 1, and then allocates tasks to processors according to their position in the list. Methods that extend on HEFT have utilised approaches such as guided search techniques to extend the schedule quality by looking ahead at decisions made by the algorithm [2]; however, these are used on top of the existing list-scheduling framework.

The literature demonstrates the efficacy of HEFT as a scheduling algorithm of good quality schedules. Since its introduction, researchers looking to develop new approaches to scheduling tasks on heterogeneous systems have built upon the algorithm with a variety of approaches, including look-ahead search and run-time scheduling [11, 2]. The combination of good schedule times and the simplicity of a list scheduling algorithm [10] has resulted in

HEFT becoming a well-established approach to scheduling on heterogeneous systems.

### 3 Methodology

#### 3.1 Representing the DAG Scheduling problem

As mentioned in Section 2, the Task Scheduling-problem can be represented formally as a graph  $G=(V,E)$ , in which  $V$  is set of  $v$  nodes and  $E$  is a set of  $e$  edges [10]. Edge weights on the graph represent communication costs between vertices, and vertex weights represent the time it takes for a task to execute. In this paper, we will use the term ‘tasks’ or ‘nodes’ interchangeably to describe vertices, following from the literature [16]. Figure 1 shows an example DAG that represents a workflow schedule, and computation costs for a heterogeneous system with two processors. The edge weights between nodes on the graph represent communication costs between nodes when transferring from one processor to another; this reflects data sizes and transfer rates from each respective task. Data transfer rates on the same processor are assumed to be negligible [9, 16, 13].

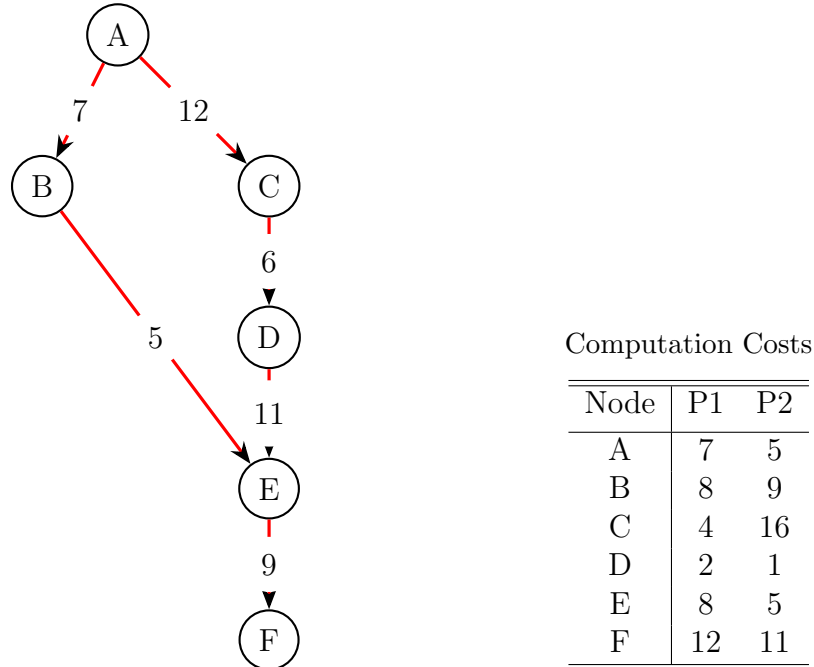


Figure 1: Example DAG and Attributes

## 3.2 Heterogeneous Earliest-Finish-Time

The HEFT algorithm is a list-scheduling algorithm that is separated into two phases: task prioritisation, and processor selection. The task prioritisation phase utilises *a priori* information about the tasks – computation, and communication costs; task precedence – to give the tasks a priority, or ‘rank’, by which they are then sorted. The processor selection phase uses this sorted list to assign tasks to processors; this is done by evaluating the Estimated Start Time (EST) of the task with respect to previously scheduled tasks (the predecessor’s Actual Finish Time (AFT)), and then assigning it to a processor that provides the earliest EST.

HEFT is different to other list-scheduling algorithms in that it not only will assign a task  $t_i$  to the earliest time on a processor  $p_j$ ; it will also iterate through the processors and check to see if  $t_i$  can be ‘inserted’ between two already scheduled tasks, provided this does not break precedence constraints. This is called the **insertion policy**. The HEFT algorithm is summarised in Algorithm 1:

---

**Algorithm 1** The HEFT algorithm

---

- 1: **procedure** HEFT\_SCHEDULE
  - 2:     Calculate mean computation and communication costs
  - 3:     Calculate upward rank for each task
  - 4:     Sort tasks by decreasing order of rank
  - 5:     **for** Each task in the sorted list of tasks **do**
  - 6:         Assign the task to a processor using the insertion policy
- 

### 3.2.1 Upward Rank

The upward rank is computed by ranking the task graph upward, from the last task (or exit node). The ranking function,  $rank_u$ , is recursively applied to each node in the graph. The function is described in Equation 1:

$$rank_u(n_i) = \overline{w_i} + \max_{n_j \in succ(n_i)} (\overline{c_{i,j}} + rank_u(n_j)) \quad (1)$$

Where  $succ(n_i)$  refers to the immediate successors of task  $n_i$ ; for example, in Figure 1, if  $n_i$  was A,  $succ(n_i)$  would be nodes B and C.  $\overline{w_i}$  and  $\overline{c_{i,j}}$  refer to the average computation and communication costs for a given node and  $(i, j)$  edge, respectively. Table 1 shows the results of applying the upward rank method to the workflow described in Figure 1.

Table 1: Results of Ranking Heuristic for DAG in Figure 1

Node	$rank_u$
A	71
B	39
C	53
D	38
E	26
F	11

It can be inferred from the table that the final sorted list of nodes does not match the cardinal order of nodes, but does retain precedence constraints; sorting by descending rank gives  $\{A, C, B, D, E, F\}$ .

### 3.2.2 Topological Sort

For any DAG, a topological sort is a list of vertices such that for every  $(i, j)$  edge,  $i$  comes before  $j$  in the list [15]. This means that a topological sort of a task DAG retains the precedence constraints of each task in the graph; given a graph with a large enough number of vertices, a topologically sorted list of nodes may not be unique, as vertices that exist at the same depth in the graph have no precedence constraints between each other. An example of this can be seen in Figure 1, where a topological sort could result in B occurring before C, or C before B.

From this definition, we can see that the sorted list provided by the  $rank_u$  calculation forms a unique topological sort for this combination of tasks and cost values.

### 3.2.3 Insertion Policy and Final Schedule

As discussed in Section 3.2, the final schedule for a given task graph is generated by inserting tasks onto processors when they have an available slot of appropriate length available. As tasks are allocated to processors based on the EST of the task, there may be a processor that is idle for a period of time. The insertion policy loops through each processor and determines if, for a given task, there is a sufficient window in which the start time and finish time of the task align. An example final schedule generated by HEFT is shown in Table 2.

Table 2: Final Schedule of DAG in Figure 1

P1	P2
B	A
	C
	D
	E
	F

We can see clearly from Figure 1 why A was scheduled to P2, as the computation cost is lower. C, the next task in the rank-sorted list is scheduled, is also scheduled on P2, even though the computation cost is greater than that on P1. However, the EST is higher on P1 due to the communication cost; in combination with the computation cost of C on P1, this would result in an AFT than that on P2. We can see that B is schedule onto P1 during the same time slot as C is on P2; given that P1 is idle, and the EST of B on P1 - including communication costs - is lower on P1 than on P2, it makes sense to schedule it there. The following tasks are all scheduled on P2, as it is trivial to see that communication costs to P1 would result in later EST times, and thus a longer schedule.

## 4 Implementation

All code used to implement and test the HEFT algorithm was written in the Python programming language, with assistance from the popular graph library NetworkX <sup>1</sup>. To assist testing during the implementation, a random DAG generator was developed to create a graph of a given number of nodes and edges; this was used in conjunction with functions that create a random computation matrix for a given number of processors, and a random communication matrix for the nodes in the graph <sup>2</sup>. In order to determine the differences between the rank and topological sorts, schedules and algorithmic execution times were collected through a loop that generated a number of different graphs with increasing number of nodes. The experimental settings were as follows:

<sup>1</sup><https://networkx.github.io>

<sup>2</sup>All code used to generate experimental data can be found at <https://github.com/myxie/heft>



- The number of nodes (N) started at 10 nodes, with an upper bound of 1000 nodes
- For each iteration of nodes, the maximum number of edges the graph could have was  $2*N$
- The step between each iteration was 50 nodes

The experiment was conducted by generating graphs in a loop, within the boundaries listed above. For each iteration, the task nodes were ranked by both the rank sort method, and a standard topological sort provided by the NetworkX library. The time data collected covered the execution time of the entire algorithm (ranking and processor selection).

## 5 Results

Results from the experiments outlined in Section 4 are found in Figure 2. Figure 2a demonstrates the difference between the final schedule length of the task graphs when using the  $rank_u$  method, or a topological sort. The plot shows a clear indication that the  $rank_u$  sorting method results in a lower schedule length than when sorted topologically. It also shows that the rate at which the length increases as the number of nodes increases is much higher for the topological sort than the ranking sort. As outlined in Section 3, there can be more than one topological sort for a given DAG. Therefore, as the number of nodes in a graph increase, the probability of the topological sort returning a less-than-ideal list of tasks also increases. This would offer an explanation for why the makespan length increases at a higher rate for the topological sort, as the rank sort is only ever generating a unique sort that has an artificial upper-bound by factoring in costs that the topological sort does not.

Figure 2b shows the difference between the execution time of the whole algorithm for a topologically sorted list of tasks, and a rank sorted list of tasks. It is clear from the plot that the  $rank_u$  sort provides a faster execution time for an increasing number of nodes. The higher execution time of the topological sort is most likely due to more time spent looping through processors during the insertion policy, as the list of nodes are not sorted based on their computation costs, and as such the time windows are not as optimally distributed between processors. Of note is the spike in execution time that occurs at 860 nodes; this is likely due to the random graph generator creating a graph that has more potential time windows to allocate as

a result of edges in the graph, as it occurs in both the topological sort and rank sort data, and the following data points continue the previous growth trend.

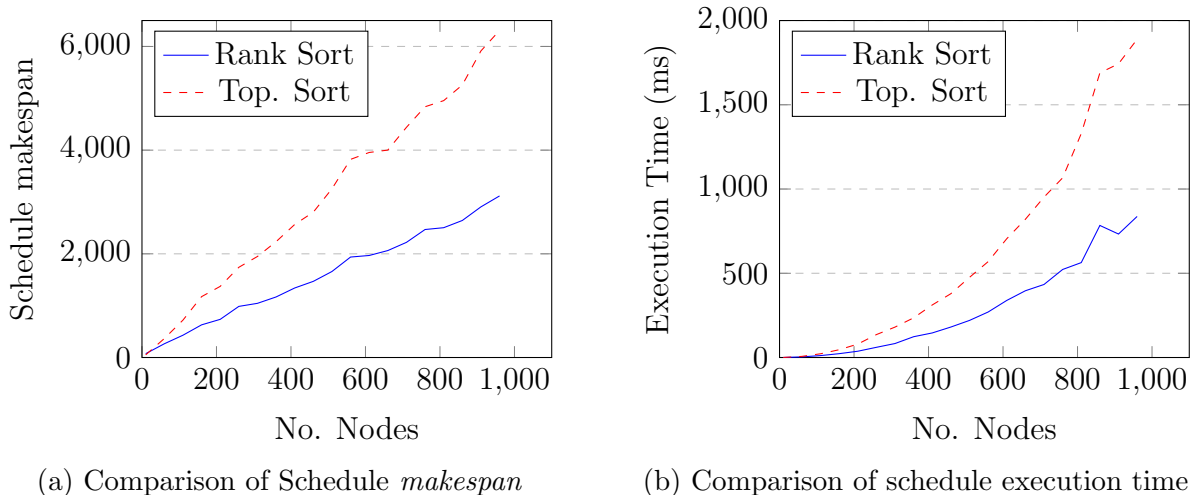


Figure 2: Comparison of Sorting methods in HEFT

## 6 Conclusion

In this paper, we presented a comparison between two ranking methods utilised in the HEFT scheduling algorithm; the original upward rank method,  $rank_u$ ; and a standard topological sort. Based on an experimental analysis run over randomly generated graphs of increasing sizes, we have shown the ranking sort provides both a better quality schedule, and a faster execution time, when compared against a topological sort. As a result, further research into improving the HEFT algorithm can focus on experimenting with and the ranking heuristic, as this paper has demonstrated it has a key influence of the overall performance of the algorithm. Planned research involves integrating the HEFT into a distributed scheduling system for large-scale scientific data processing [19]; understanding that the impact the quality of the task-ranks has on the final schedule of the system is pivotal in ensuring the appropriate information is gathered for scheduling in a high-performance computing environment.

## References

- [1] ARABNEJAD, H., AND BARBOSA, J. G. List scheduling algorithm for heterogeneous systems by an optimistic cost table. *IEEE Transactions on Parallel and Distributed Systems* 25, 3 (March 2014), 682–694.
- [2] BITTENCOURT, L. F., SAKELLARIOU, R., AND MADEIRA, E. R. M. Dag scheduling using a lookahead variant of the heterogeneous earliest finish time algorithm. In *2010 18th Euromicro Conference on Parallel, Distributed and Network-based Processing* (Feb 2010), pp. 27–34.
- [3] BLAZEWICZ, J., LENSTRA, J., AND KAN, A. Scheduling subject to resource constraints: classification and complexity. *Discrete Applied Mathematics* 5, 1 (1983), 11 – 24.
- [4] CANON, LOUIS-CLAUDE AND JEANNOT, EMMANUEL AND SAKELLARIOU, RIZOS AND ZHENG, WEI”, EDITOR=”GORLATCH, SERGEI AND FRAGOPOULOU, PARASKEVI AND PRIOL, THIERRY. *Comparative Evaluation Of The Robustness Of DAG Scheduling Heuristics*. Springer US, Boston, MA, 2008, pp. 73–84.
- [5] CASAVANT, T. L., AND KUHL, J. G. A taxonomy of scheduling in general-purpose distributed computing systems. *IEEE Transactions on Software Engineering* 14, 2 (Feb 1988), 141–154.
- [6] CHAUDHARY, V., AND AGGARWAL, J. K. A generalized scheme for mapping parallel algorithms. *IEEE Transactions on Parallel and Distributed Systems* 4, 3 (Mar 1993), 328–346.
- [7] HWANG, R., GEN, M., AND KATAYAMA, H. A comparison of multiprocessor task scheduling algorithms with communication costs. *Computers Operations Research* 35, 3 (2008), 976 – 993. Part Special Issue: New Trends in Locational Analysis.
- [8] JING-CHIOU LIOU AND M. A. PALIS. A Comparison of General Approaches to Multiprocessor Scheduling. In *Proceedings 11th International Parallel Processing Symposium* (Apr 1997), pp. 152–156.
- [9] KWOK, Y.-K., AND AHMAD, I. Benchmarking and comparison of the task graph scheduling algorithms. *Journal of Parallel and Distributed Computing* 59, 3 (1999), 381 – 422.
- [10] KWOK, Y.-K., AND AHMAD, I. Static scheduling algorithms for allocating directed task graphs to multiprocessors. *ACM Comput. Surv.* 31, 4 (Dec. 1999), 406–471.

- [11] LEE, L. T., CHANG, H. Y., LIU, K. Y., CHANG, G. M., AND LIEN, C. C. A dynamic scheduling algorithm in heterogeneous computing environments. In *2006 International Symposium on Communications and Information Technologies* (Oct 2006), pp. 313–318.
- [12] MICHELL, G. D., AND GUPTA, R. K. Hardware/software co-design. *Proceedings of the IEEE* 85, 3 (Mar 1997), 349–365.
- [13] NIKHIL, R., ET AL. Executing a program on the mit tagged-token dataflow architecture. *IEEE Transactions on computers* 39, 3 (1990), 300–318.
- [14] SARKAR, V. *Partitioning and Scheduling Parallel Programs for Multiprocessors*. Pitman Publishing, London, UK, 1989.
- [15] SKIENA, S. S. *The Algorithm Design Manual: Text*, vol. 1. Springer Science & Business Media, 1998.
- [16] TOPCUOGLU, H., HARIRI, S., AND WU, M.-Y. Performance-effective and low-complexity task scheduling for heterogeneous computing. *IEEE Transactions on Parallel and Distributed Systems* 13, 3 (Mar 2002), 260–274.
- [17] ULLMAN, J. D. Np-complete scheduling problems. *J. Comput. Syst. Sci.* 10, 3 (June 1975), 384–393.
- [18] VIJAYALAKSHMI, S R; PADMAVATHI, G. A performance study of ga and lsh in multiprocessor job scheduling. *International Journal of Computer Science Issues* 7, 1 (2010), 37–42.
- [19] WU, C., TOBAR, R., VINSEN, K., WICENEC, A., PALLOT, D., LAO, B., WANG, R., AN, T., BOULTON, M., COOPER, I., DODSON, R., DOLENSKY, M., MEI, Y., AND WANG, F. DALiuGE: A Graph Execution Framework for Harnessing the Astronomical Data Deluge. *ArXiv e-prints* (Feb. 2017).