

Review of Graph Based Scheduling Algorithms

Gagandeep Singh¹, Gurjyot Kaur², Dr.Gurvinder Singh³

¹Research Scholar, Department of Computer Science, Guru Nanak Dev University, Amritsar

³Head of Department, Department of Computer Science,, Guru Nanak Dev University, Amritsar

Abstract:For the problem of task scheduling in a parallel program the weighted directed acyclic graph(DAG) is used with a set of homogenous processors for the completion of program in the minimized time. There has been a great demand for high speed computing in many application areas. High speed electronic devices were developing to achieve high computing speed. But, the processing speed of uniprocessor computer can't be increased beyond a limit (few million floating point operations per second) because of the physical limitations imposed by the electrical properties of electronics devices [1]. To achieve high performance, major developments and improvements in the field of processing techniques and computer architectures have been made. Thus, parallel processing approach had gradually emerged to meet the computational requirement of various problem and also to enhance the efficiency of solving various current applications like weather forecasting, military defense, medical diagnosis, simulation, etc.

There are many researchers have proposed so many algorithms for scheduling problems, most of these are reported to be efficient, but it is not clear how they compare against each other. Due to many number of issues, it is difficult to measure the actual performance of algorithms and then comparison of these algorithms is also a difficult task. First, most of the scheduling algorithms are based upon diverse assumptions, making the performance comparison rather purposeless. Second, standard set of benchmarks to examine these algorithms is not available due to which the actual performance evaluation is not possible. Third, small problem size cannot evaluate the actual performance of algorithms. The main purpose of this paper is to provide the taxonomy for classifying various algorithms of different categories according to their assumptions and functionalities.

In this paper we discussed 14 scheduling algorithms and give the characteristics of all algorithms.

Keywords: Multiprocessors, parallel processing, scheduling, task graphs, scalability.

1. INTRODUCTION

The problem of scheduling a DAG, also called a taskgraph or macro-dataflow graph. In the DAG a set of nodes is used. The set of nodes are homogeneous processors, are used for scheduling the tasks for execution in the minimized time [11]. The scheduling problem is NPcomplete in its general forms [9], and only for a few restricted cases the polynomial time solutions are available [4], [8]. It is difficult to schedule the problem in an efficient manner for achieving a meaningful speedup in a parallel or distributed system, so it is important to make the interest of researchers of all research community to make the efficient execution algorithms. Considerable research effort expended in solving the problem by using many heuristic algorithms. While each heuristic is individually reported to be efficient, it is not clear how these algorithms compare against each other on a unified basis.

The objectives of this study include the 14 different algorithms under the categories of UNC, BNP, APN and their characteristics. In this literature survey the large number of DSAs have been presented with different assumptions, it is important to define these algorithms into various classes according to their assumptions about the program and machine model.

1.1 DAG Model

The DAG is a generic model of a parallel program consisting of a set of processes among which there are dependencies. Each process is expressed by an atomic node. An atomic node has one or more inputs. When all inputs are available, the node is triggered to execute. After its execution, it generates its outputs. In this model, a set of v nodes $\{n_1, n_2, \dots, n_n\}$ are connected by a set of e directed edges, each of which is denoted by (n_i, n_j) , where n_i is called the parent and n_j is called the child. A node without parent is called an entry node and a node without child is called an exit node. The weight of a node, denoted by $w(n_i)$, is equal to the process execution time. Since each edge corresponds to a message transfer from one process to another, the weight of an edge, denoted by c

(n_i, n_j) is equal to message transmission time. Thus, $c(n_i, n_j)$ becomes zero when n_i and n_j are scheduled to the same processor communication time[2].

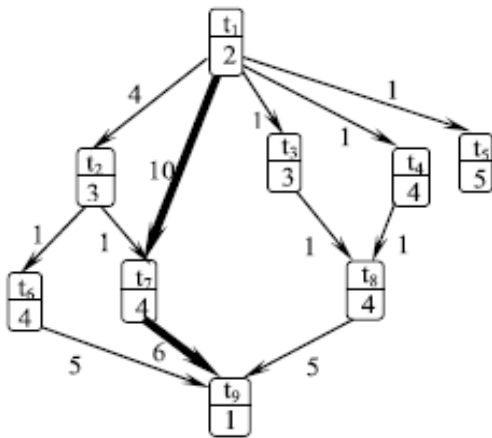


Fig.1 Directed acyclic graph

Consider above directed acyclic task graph $G = \{V, E\}$ of n nodes. Each node $V = \{T_1, T_2, \dots, T_n\}$ in the graph represents a task. Aim is to map every task to a set $P = \{P_1, P_2, \dots, P_m\}$ of m processor. Each task T_i has a weight W_i associated with it, which is the amount of time the task takes to execute on any one of the m homogeneous processors. Each directed edge e_{ij} indicates dependence between the two tasks T_i and T_j that it connects. If there is a path from node T_i to node T_j in the graph G , then T_i is the predecessor of T_j and T_j is the successor of T_i . The successor task cannot be executed before all its predecessors have been executed and their results are available at the processor at which the successor is scheduled to execute. A task is "ready" to execute on a processor if all of its predecessors have completed execution and their results are available at the processor on which the task is scheduled to execute. If the next to be executed on a processor is not yet ready, the processor remains idle until the task is ready. The elements set C are the weights of the edges $s C = \{c_k; k = 1, 2, 3 \dots r\}$ It represents the data communication between the two tasks, If they are scheduled to different processors. But if both tasks are schedule to the same processor, then the weight associated to the edge becomes null[13].

If the node n_i scheduled on the processor P then the start time and the finish time of node n_i is denoted by $ST(n_i, P)$ and $FT(n_i, P)$ respectively. The main objective of the DAG scheduling is to find the earlier start time of the tasks to the processors such that schedule length minimized such that the precedence constraints are preserved.

For solving the problem of DAG scheduling mostly the researchers are interested to design efficient heuristics which can find the better solution for the problem within a reasonable amount of time. Most of these heuristic algorithms are based on the *list scheduling* techniques. The concept of priority assignment is used. The priorities are assigned to the nodes and nodes are arranged in descending order of priorities. Higher priority node is scheduled before lower priority node.

Assigning Priorities to Nodes: The *t-level* (top level) and *b-level* (bottom level) are the two major attributes which are used to assigning priorities. The *t-level* of a node n_i is the length of the longest path from an entry node n_i to in the DAG (excluding n_i). Here, the length of a path is the sum of all the node and edge weights along the path. The *t-level* of n_i highly correlates with n_i 's earliest start time, denoted by $TS(n_i)$, which is determined after n_i is scheduled to a processor.

node	SL	t-level	b-level	ALAP
*n ₁	11	0	23	0
n ₂	8	6	15	8
n ₃	8	3	14	9
n ₄	9	3	15	8
n ₅	5	3	5	18
n ₆	5	10	10	13
*n ₇	5	12	11	12
n ₈	5	8	10	13
*n ₉	1	22	1	22

Fig.2

(a) The Static Levels (SLs) t-levels, b-levels and ALAP

The *b-level* of a node is the length of the longest path from node to an exit node and is bounded by the length of the *critical path*. The path from an entry node to an exit node, whose length is the maximum is called *critical path* (CP) of a DAG.

In the *b-level* of a node variations in the computations are possible. Most DSAs examine a node for scheduling only after all the parents of the node have been scheduled. In this case, the *b-level* of a node is a constant until after it is scheduled to a processor. However, some algorithms allow the scheduling of a child before its parents. In that case, the *b-level* of a node becomes a dynamic attribute. Different DSAs have used the *t-level* and *b-level* attributes in a variety of ways. Some algorithms assign a higher priority to a node with a larger *b-level* while some algorithms assign a higher priority to a node with a smaller *t-level*.

For determining the start time of a node on a processor P , some algorithms only consider scheduling a node after the last node on P . Some algorithms also consider other idle time slots on P and may insert a node between two already scheduled nodes.

Critical-Path-Based vs. Non-Critical-Path-Based:

Critical-path-based algorithms give a higher priority to a critical-path node (CPN). Noncritical-path-based algorithms simply assign priorities based on the levels of the nodes.

Static List vs. Dynamic List: The ready list maintained a set of ready nodes. The ready list is sorted in descending order. Initially, only the entry nodes are included in the ready list. After a node is scheduled, the next node which is not scheduled is inserted into the ready list.

The list can be maintained in two ways: static ready list and dynamic ready list. In the static ready list the list is constructed before scheduling starts and remains the same throughout the whole scheduling process and in the dynamic ready list the list is rearranged according to the changing node priorities.

Greedy vs. Non-Greedy: Most scheduling algorithms attempt to minimize the start-time of a node for assigning a node to a processor. This is a greedy strategy. But in Non greedy, the algorithms do not minimize the start-time of a node but consider other factors as well.

Time-Complexity: The some terms like number of node, the number of edges, and the number of processors are used to calculate the time-complexity of a DSA. The main step of algorithm is to traverse the DAG and search of slots in the processors to place a node. The dynamic priority assignment having higher time-complexity as compared to static priority assignment.

2. LITERATURE SURVEY

(Adam *et al*, 1974) discussed the comparison of list scheduling algorithms for the parallel processing system. This paper describe that how list scheduling algorithms vary their performance in the different situations in the parallel environment.

(Ahmed *et al*, 1997) is describe the concept of automatic parallelization an scheduling of programs on multiprocessor environment using CASCH.

(Grahmet *al*, 1972) is discussed the optimal ways of scheduling of tasks on Two-Processor system.

(Loiet *al*, 1995) is discussed the techniques of automatic graph generation for scheduling the tasks on the processors in the efficient manner so that the total time of execution can be minimized.

(Lewis *et al*, 1990) this paper discussed the concept of scheduling of parallel programs on the arbitrary number of

machines. This paper also gives the comparison analysis of the performance of the scheduling when the number of processors are non uniform.

(Fernandez *et al*, 1973) discussed the concept of BNP scheduling. It also describes time for multiprocessors optimal schedule for the bounded number of processors.

(Hwang *et al*, 1989) discussed the Scheduling of the precedence graph having the inter-processor communication time in the multiprocessor environment.

(Kim *et al*, 1988) This paper describes the General approach to mapping of parallel computation upon multiprocessor architectures.

(Gurvinder Singh *et al*, 2011) It is the survey paper about all the scheduling algorithms that allocates parallel program to DAG on homogenous processors. The main objective is to minimize the execution time, evaluation of performances and comparison of all algorithms based on their performances. The BNP, TDB, UNP algorithms are compared in this paper.

(Sachiet *al*, 2013) discussed a literature in which several heuristic methods have been developed that obtain suboptimal solutions in less than the polynomial time.

3. TAXONOMY OF DAG SCHEDULING ALGORITHMS

The list scheduling algorithms are divided into two categories. Some of these algorithms assume the uniform computational costs of all the tasks ([5], [20]) whereas some of these algorithms assume the arbitrary computational costs. The earlier work done by researchers the inter-task communication assumed to be zero, that is, the task graph contains precedence but without cost. The problem becomes less complex in the absence of communication delays.

For minimizing the communication delays of scheduling the task duplication is used. Duplication is done by duplicating the ancestor nodes on which the predecessors of those nodes are dependent for the execution.

There are two former classes of algorithms of list scheduling problems are called the UNC i.e. *unbounded number of clusters* scheduling algorithms [2] and the BNP i.e. *bounded number of processors* scheduling algorithms [2]. In both classes of algorithms, the processors are assumed

to be fully-connected and no attention is paid to link contention or routing strategies used for communication.

3.1. UNC Scheduling Algorithm

The concept of clustering is used in UNC algorithms. Each node in UNC considered as a cluster. In this algorithm the merging of the clusters is done if it reduces the completion time. Merging continuously performed until no cluster can be merged. The idea behind UNC algorithm is that it uses more processors to further reducing the schedule length. The post processing step is performed for mapping the clusters onto the processors because the cluster may be more than the number of processors.

In UNC scheduling we have study the five algorithms named as EZ, LC, DSC, MD, DCP and study their characteristics.

The EZ Algorithm: This Edge Zeroing algorithm [15] based on edge weights clusters are selected and then merging those clusters. The algorithm finds the edge with the largest weight at each step. The two clusters incident by the edge are merged and the merging is done if and only if it does not increase the completion time. After merging the two clusters, the SLs of the nodes are used for ordering of the nodes in the resultant cluster.

The LC Algorithm: This Linear Clustering algorithm [13] forming a single cluster based on the CP by merging the nodes. In the first step the set of nodes constituting the CP is determined by the algorithm. Then it schedules all of the CP nodes to a single processor at once. The scheduled nodes and all incident edges are then removed from the DAG. The algorithm zeroes the edges on the entire CP at once.

The CP may change when an edge is zeroed. The original CP is not containing the edge that should be zeroed next.

The DSC Algorithm: The Dominant Sequence Clustering algorithm [9] considers the Dominant Sequence of a graph. It is simply the CP of the partially scheduled DAG. The DSC algorithm tracks the CP of the partially scheduled DAG at

each step. For this it uses the composite attribute (*b-level + t-level*) as the priority of a node. Unless the node is ready, the DSC algorithm does not select the node having the highest priority for scheduling.

The MD Algorithm: If a node is on the current CP of the partially scheduled DAG, the sum of its *b-level* and *t-level* is equal to the current CP length. Thus, the relative mobility of a node is zero if it is on the current CP. At each step, the MD algorithm selects the node with the smallest relative mobility for scheduling. In testing whether a cluster can accommodate a node, the MD algorithm scans from the earliest idle time slot on the cluster and schedules the node into the *first* idle time slot that is large enough for the node.

The DCP Algorithm: The Dynamic Critical Path algorithm [14] is designed based on a similar attribute to relative mobility. To find a better cluster this algorithm uses a look-ahead strategy

for a given node. In addition, The DCP algorithm computes the value of $T_S(n_c)$ on a cluster for computing the value of $T_S(n_i)$ on the same cluster, where, n_i is the child of n_i has the largest communication and is called the critical child of n_i . The DCP algorithm schedules n_i to the cluster that gives the minimum value of the sum of these two attributes. This look-ahead strategy can potentially avoid scheduling a node to a cluster

that has no room to accommodate a heavily communicated child of the node. The DCP algorithm examines all the existing clusters for a node while the MD algorithm only tests from the first cluster and stops after finding one that has a large enough idle time slot.

Algorithm	Priority	List	CP Based	Greedy	Complexity
EZ	SL	Dynamic	No	No	$O(e(v + e))$
LC	SL+ t-level	Static	Yes	No	$O(v(v + e))$
DSC	SL+ t-level	Dynamic	Yes	Yes	$O((e + v)\log v)$
MD	b-level+t-level	Dynamic	Yes	No	$O(v^3)$
DCP	b-level+t-level	Dynamic	Yes	No	$O(v^3)$

Table 1: UNC Scheduling Algorithm and their characteristics

are given in the table below where p denotes the number of processors given.

3.2. BNP Scheduling Algorithms[3]

In this scheduling technique we have discussed five algorithms: HLET, MCP, ISH, ETF, LAST. The characteristics

The HLFET Algorithm: The HLFET is the simplest scheduling algorithms. In this algorithm the processor that allows the earliest start time is used for scheduling a node.

The main drawback of HLFET is that it calculating the SL of a node and ignores the communication costs on the edges.

The ISH Algorithm: The ISH i.e Insertion Scheduling Heuristic algorithm [14] uses the holes created by the partial schedules with simple but effective idea. It first picks an unscheduled node. The processor that allows the earliest start time is used to schedule the unscheduled node having highest SL. It tries to insert other unscheduled nodes from the ready list into the idle time slot before the node just scheduled.

The MCP Algorithm: The MCP i.e Modified Critical Path algorithm [16] uses the ALAP time of a node as a priority. For computing the ALAP time of a node, first computing the length of CP and then subtracting the *b-level* of the node from it. Thus, *t-level* on the CP are the ALAP times of the nodes. After computing the ALAP times of all the nodes, MCP algorithm then constructs a list in ascending order of ALAP times of all the nodes.

Ties are broken by considering the ALAP times of the children of a node. Then one by one the algorithm schedules the nodes on the list such that the insertion approach is used for scheduling and the processor that allows the earliest start time is used for scheduling a node.

Algorithms	Priority	CP-Based	List Type	Greedy	Complexity
HLEFT	SL	No	Static	Yes	$O(v^2)$
ISH	SL	No	Static	Yes	$O(v^2)$
MCP	ALAP	Yes	Static	Yes	$O(v^2 \log v)$
ETF	SL	No	Static	Yes	$O(pv^3)$
LAST	Edge weights	No	Dynamic	Yes	$O(v(v + e))$

Table 2: BNP Scheduling Algorithm and their characteristics

3.3. APN Scheduling Algorithms

APN Scheduling we have discussed four algorithms MH, DLS, BU, BSA and the characteristics are given in the below table.

The MH Algorithm: The MH i.e Mapping Heuristic algorithm [7] initially makes a ready node list. It contains entries of all nodes arranged in descending order according to their priorities. The processor having the smallest start time is used for scheduling the node. A routing table is used in this algorithm which maintained the calculated start time of a node for each processor. The table contains the information from the parent nodes to the nodes under consideration. When a node is scheduled on processor all of its ready successor nodes are appended to the ready node list.

The ETF Algorithm: The ETF i.e Earliest Time First algorithm [10] selects the node having smallest start time. For finding smallest start time firstly earliest start time of all the nodes are computed at each step. The tie of two nodes having same earliest start time breaks by scheduling the one with the higher SL. Thus, a node with a higher SL does not necessarily get scheduled first because node with the earliest start time having higher priority according to the algorithm.

The LAST Algorithm: The LAST i.e Localized Allocation of Static Tasks algorithm [9] is not a list scheduling algorithm. It uses an attribute called D_NODE for the node priority. It depends on the incident edges of a node. The main Objective of this algorithm is to reduce the overall communication. In this algorithm the node can be selected before its parent for scheduling. One of the consequences of using D_NODE is that a node may be selected before some of its parents for the scheduling. Thus, until the scheduling process terminates the earliest start time of a node cannot be fixed. For the node selection process the node weight is ignored in the LAST algorithm.

The DLS Algorithm: The DLS i.e Dynamic Level Scheduling algorithm [11] is used as an APN scheduling algorithm. In the APN scheduling algorithm, The message routing method are required for DLS. These routing methods are supplied by the user. Then based on that how the message is routed from the parents of the node, the T_s of a node is computed

The BU Algorithm: The BU i.e Bottom-Up algorithm [12] for assigning the nodes on the CP it first finds the CP of DAG, and then assigns all the nodes of CP on the same processor at once. The reversed topological order of nodes is used for assigning the remaining nodes to the processors. The load balancing of processors are required at the time of assignment of nodes. The processor selection is based on some heuristics to balance the load across all given processors. The BU algorithm tries to schedule the

communication messages among all the nodes assigned to processors using a channel allocation heuristic to keep the hop count of every message. Different network topologies require different channel allocation heuristics.

The BSA Algorithm: The BSA Bubble Scheduling and Allocation algorithm constructs a schedule incrementally. For this it first injecting all the nodes to the pivot processor, defined as the processor with the highest degree. Then algorithm tries to improve the start time of each node by transferring it to one of the

adjacent processor of the pivot processor if the migration can improve the start time of the node. This is because after a node migrates, the space it occupies on the pivot processor is released and can be used for its successor nodes on the pivot processor. After all nodes on the pivot processor are considered, the algorithm selects the next processor in the processor list to be the new pivot processor. The process is repeated by changing the pivot processor in a breadth-first order.

Algorithms	Priority	CP-Based	Message Routing	Complexity
MH		SL	No	$O(v(p^3v + e))$
DLS		SL-T _s	Yes	$O(v^3p^2)$
BU		-----	Yes	$O(v^2 \log v)$
BSA		-----	Yes	$O(p^2ev)$

Table 3: APN Scheduling Algorithm and their characteristics

4. CONCLUSIONS AND FUTURE WORK

In this paper we have studied 14 different algorithms for the DAG scheduling problems and their characteristics are discussed. After this we have concluded that:

- As demonstrated by both DCP and DSC algorithm the dynamic critical path is better than the static critical path.
- In general the dynamic priority is better than static priority.
- As compared to both BNP and UNC algorithms APN algorithm is complicated because it takes more parameters.
- In UNC algorithms the bounded numbers of processors are used to assigning the clusters obtained through scheduling and all nodes of cluster assigned to a same processor. Due to this property cluster scheduling algorithms become more complex as compared to BNP scheduling algorithms.
- BNP and UNC classes' algorithms having more accurate performance than other scheduling algorithms.

In the future work these algorithms can be implemented using the Genetic Approach (GA). The Genetic algorithms (Gas) are adaptive heuristic search based algorithms based on the evolutionary idea of natural selection and genetics. GAs are a part of evolutionary computing, a rapidly growing area of artificial intelligence. The scheduling algorithms with Genetic Approach give better performance than the simple scheduling algorithms.

5. REFERENCES

- [1] T.L. Adam, K.M. Chandy, and J. Dickson, "A Comparison of List Scheduling for Parallel Processing Systems," *Comm. ACM*, vol. 17, no. 12, pp. 685-690, Dec. 1974.
- [2] I. Ahmad, Y.-K. Kwok, M.-Y. Wu, and W. Shu, "Automatic Parallelization and Scheduling of Programs on Multiprocessors using CASCH," *Proc. 1997 Int'l Conf. Parallel Processing*, pp. 288-291, Aug. 1997.
- [3] Parneet Kaur, Dheerendra Singh, Gurvinder Singh and Navneet Singh, "Analysis, comparison and performance evaluation of BNP scheduling algorithms in parallel processing", *International Journal of Information Technology and Knowledge Management* January-June 2011, Volume 4, No. 1, pp. 279-284,
- [4] Sachi Gupta, Gaurav Agarwal, and Vikas Kumar, "An Efficient and Robust Genetic Algorithm for Multiprocessor Task Scheduling", *International Journal of Computer Theory and Engineering*, Vol. 5, No. 2, April 2013.
- [5] E.G. Coffman and R.L. Graham, "Optimal Scheduling for Two-Processor Systems," *Acta Informatica*, vol. 1, pp. 200-213, 1972.
- [6] M. Cosnard and M. Loi, "Automatic Task Graph Generation Techniques," *Parallel Processing Letters*, vol. 5, no. 4, pp. 527-538, Dec. 1995.
- [7] H. El-Rewini and T.G. Lewis, "Scheduling Parallel Programs onto Arbitrary Target Machines," *J. Parallel and Distributed Computing*, vol. 9, no. 2, pp. 138-153, June 1990.
- [8] E.B. Fernandez and B. Bussell, "Bounds on the Number of Processors and Time for Multiprocessor Optimal Schedules," *IEEE Trans. Computers*, vol. C-22, no. 8, pp. 745-751, Aug. 1973.

- [9] M.R. Garey and D.S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W.H. Freeman and Company, 1979.
- [10] J.J. Hwang, Y.C. Chow, F.D. Anger, and C.Y. Lee, "Scheduling Precedence Graphs in Systems with Interprocessor Communication Times," *SIAM J. Computing*, vol. 18, no. 2, pp. 244-257, Apr. 1989.
- [11] K. Hwang and Z. Xu, *Scalable Parallel Computing: Technology, Architecture, Programming*, McGraw-Hill, 1998.
- [12] Y.C. Chung and S. Ranka, "Application and Performance Analysis of a Compile-Time Optimization Approach for List Scheduling Algorithms on Distributed-Memory Multiprocessors," *Proc. Supercomputing'92*, pp. 512-521, Nov. 1992.
- [13] S.J. Kim and J.C. Browne, "A General Approach to Mapping of Parallel Computation upon Multiprocessor Architectures," *Proc. 1988 Int'l Conference on Parallel Processing*, vol. II, pp. 1-8, Aug. 1988.
- [14] B. Kruatrachue and T.G. Lewis, "Duplication Scheduling Heuristics (DSH): A New Precedence Task Scheduler for Parallel Processor Systems," Technical Report, Oregon State University, Corvallis, OR 97331, 1987.
- [15] V. Sarkar, *Partitioning and Scheduling Parallel Programs for Multiprocessors*, MIT Press, Cambridge, MA, 1989.
- [16] M.-Y. Wu and D.D. Gajski, "Hypertool: A Programming Aid for Message-Passing Systems," *IEEE Trans. Parallel and Distributed Systems*, vol. 1, no. 3, pp. 330-343, July 1990.