

Performance-Effective and Low-Complexity Task Scheduling for Heterogeneous Computing

Haluk Topcuoglu, Salim Hariri, Min-You Wu

Jay Khatri Keshav Parihar Prateekshya Priyadarshini
Vijay Purohit G Vasantha Reddy

INDIAN INSTITUTE OF TECHNOLOGY, GUWAHATI

April 21, 2022



INTRODUCTION

Heterogeneous Computing

Diverse sets of resources interconnected with high speed network.

General Task Scheduling Problem

To assign tasks to suitable processors and then order execution of the tasks.

Static Task Scheduling Problem

- ◊ Applications → DAGs
- ◊ Application Tasks → Nodes
- ◊ Inter-task data dependencies → Edges
- ◊ Computation costs → Node labels
- ◊ Communication costs → Edge labels
- ❑ NP-Complete in the general/restricted cases.

TASK SCHEDULING PROBLEM

Scheduling System Model

Application Target Computing Environment Performance Criteria

- ◊ $DAG(V, E)$, V : vertices, E : edges, $|V| = v$ tasks, $|E| = e$
- ◊ $(i, j) \in E$ means n_i should finish its execution before n_j starts
- ◊ Data: $v \times v$ Matrix, $data_{i,j}$: amount of data to be transmitted from n_i to n_j
- ◊ Entry Task: A task without any parent.
- ◊ Exit Task: A task without any child.
- ◊ Zero-Cost pseudo exit (entry) task.



TASK SCHEDULING PROBLEM

Scheduling System Model

Application Target Computing Environment Performance Criteria

- ◊ Q : set of heterogeneous processors, $|Q| = q$.
 - ▶ Connected in fully connected topology.
 - ▶ Inter-processor communications without contention.
 - ▶ Assumption: computation overlap with communication.
 - ▶ Non-preemptive task execution.
- ◊ W : $v \times q$ matrix, $w_{i,j}$: estimated execution time to complete n_i on p_j
- ◊ B : $q \times q$ matrix, $B_{i,j}$: data transfer rates between p_i and p_j
- ◊ L : q dimensional vector, $q(i)$: communication startup cost for p_i



TASK SCHEDULING PROBLEM

Scheduling System Model

Application Target Computing Environment Performance Criteria

Objective Function

Assign tasks to suitable processors and order their executions maintaining the precedence requirements and minimizing the overall completion time.



TASK SCHEDULING PROBLEM

TARGET COMPUTING ENVIRONMENT

Average Execution Cost $\overline{w_i}$

$$\overline{w_i} = \sum_{j=1}^q \frac{w_{i,j}}{q}$$

- $\overline{w_i}$: average execution cost of n_i
- $w_{i,j}$: estimated execution time to complete n_i on p_j
- q : number of processors



TASK SCHEDULING PROBLEM

TARGET COMPUTING ENVIRONMENT

Communication Cost

$$c_{i,k} = L_m + \frac{\text{data}_{i,k}}{B_{m,n}}$$

- $c_{i,k}$: communication cost of the edge (i, k) for transferring data from n_i (on p_m) to n_k (on p_n)
- L_m : communication startup cost of p_m
- $\text{data}_{i,k}$: amount of data to be transmitted from n_i to n_k
- $B_{m,n}$: data transfer rate between P_m and P_n



TASK SCHEDULING PROBLEM

TARGET COMPUTING ENVIRONMENT

Average Communication Cost

$$\overline{c_{i,k}} = \overline{L} + \frac{\text{data}_{i,k}}{\overline{B}}$$

- $\overline{c_{i,k}}$: average communication cost of edge (i, k)
- \overline{L} : average communication startup time
- $\text{data}_{i,k}$: amount of data to be transmitted from n_i to n_k
- \overline{B} : average transfer rate among the processors



EARLIEST EXECUTION START TIME

Earliest Execution Start Time

$$\text{EST}(n_{\text{entry}}, p_j) = 0$$

$$\text{EST}(n_i, p_j) = \max \left\{ \text{avail}[j], \max_{n_m \in \text{pred}(n_i)} (\text{AFT}(n_m) + c_{m,i}) \right\}$$

- $\text{EST}(n_i, p_j)$: earliest execution start time of n_i on p_j
- $\text{avail}[j]$: earliest time at which p_j is ready for task execution
- $c_{m,i}$: communication cost of the edge (m, i) for transferring data from n_m to n_i
- $\text{AFT}(n_m)$: actual finish time of n_m
- $\text{pred}(n_i)$: set of immediate predecessor tasks of n_i



TASK SCHEDULING PROBLEM

TARGET COMPUTING ENVIRONMENT

Earliest Execution Finish Time

$$EFT(n_i, p_j) = w_{i,j} + EST(n_i, p_j)$$

- $EFT(n_i, p_j)$: earliest execution finish time of n_i on p_j
- $w_{i,j}$: estimated execution time to complete n_i on p_j
- $EST(n_i, p_j)$: earliest execution start time of n_i on p_j



TASK SCHEDULING PROBLEM

TARGET COMPUTING ENVIRONMENT

Schedule Length (Makespan)

$$\text{makespan} = \max \{ \text{AFT}(n_{\text{exit}}) \}$$

- *makespan*: schedule length (overall completion time)
- *AFT*(n_{exit}): actual finish time of n_{exit}



RELATED WORK

FOR HOMOGENEOUS PROCESSORS

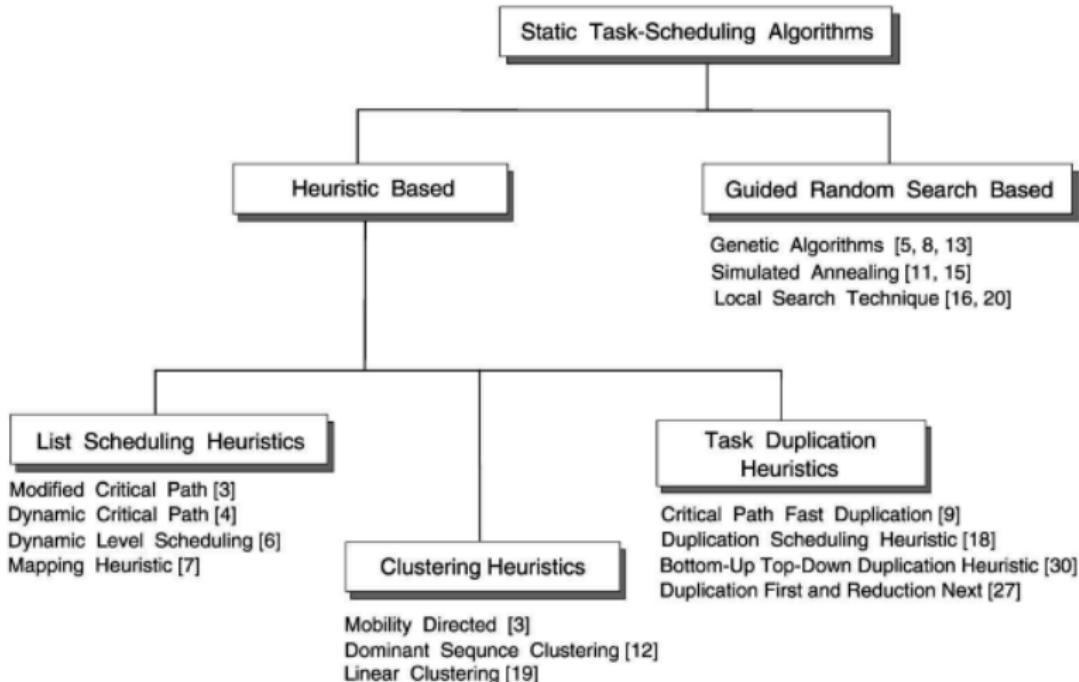


Figure: Classification of Static Task-Scheduling Algorithms.



STATIC TASK-SCHEDULING ALGORITHMS

RELATED WORK FOR HOMOGENEOUS PROCESSORS

List Scheduling Heuristics

- Ordered list of tasks by assigning priority.
- Two Phases: Task Prioritizing and Processor Selection
- Selection based on priority
- Schedule on processor to minimize predefined cost function.
- Examples: Dynamic Critical Path (DCP), Mapping Heuristic (MH).
- Advantages:
 - ▶ Comparable with other categories.
 - ▶ More Practical.
- Disadvantages:
 - ▶ For bounded number of fully connected homogeneous processors.



STATIC TASK-SCHEDULING ALGORITHMS

RELATED WORK FOR HOMOGENEOUS PROCESSORS

Clustering Heuristics

- Steps: Cluster Merging, Cluster Mapping, Task Ordering Step.
- Same Cluster Task executed on same processor.
- Examples: Clustering and Scheduling System (CASS), Dominant Sequence Clustering (DSC).
- Disadvantages:
 - ▶ Requires Additional Steps.
 - ▶ For unbounded number of processors. Not directly applicable.



STATIC TASK-SCHEDULING ALGORITHMS

RELATED WORK FOR HOMOGENEOUS PROCESSORS

Task Duplication-Based Heuristics

- Mapping of some tasks redundantly.
- Examples: BTDH and DSH Algorithm $O(v^4)$.
- Disadvantages:
 - ▶ High Time Complexity.
 - ▶ For unbounded number of identical processors.



STATIC TASK-SCHEDULING ALGORITHMS

RELATED WORK FOR HOMOGENEOUS PROCESSORS

Guided Random Search Techniques

- Use random choice to guide themselves through the problem space.
- Combine the knowledge gained from previous search results with randomizing features.
- Examples: Genetic Algorithms (GAs), Simulated Annealing.
- Advantages:
 - ▶ Popular and widely studied used techniques.
 - ▶ generate good quality of output schedules.
- Disadvantages:
 - ▶ High scheduling time.
 - ▶ Control Parameters should be determined appropriately.



RELATED WORK

Static Task-Scheduling Algorithms For Heterogeneous Environments

Dynamic-Level Scheduling (DLS) Algorithm

- (ready node, available processor) pair that maximizes the dynamic level value.
- Dynamic Value

$$DL(n_i, p_j) = rank_u^s(n_i) - EST(n_i, p_j)$$

- Computation cost of a task is the median value of the computation costs of the task on the processors.
- Disadvantage:
 - ▶ Time Complexity $O(v^3 \times q)$



RELATED WORK

Static Task-Scheduling Algorithms For Heterogeneous Environments

Mapping Heuristic (MH)

- computation cost of a task on processor = num of instructions to be executed/speed of the processor.
- Assumed homogeneous processor before scheduling.
- Uses static upward ranks to assign priorities.
- Ready time of a processor for a task is the time when the processor has finished its last assigned task and is ready to execute a new one.
- Disadvantages:
 - Time Complexity $O(v^2 \times q^3)$ with contention.



RELATED WORK

Static Task-Scheduling Algorithms For Heterogeneous Environments

Levelized-Min Time (LMT) Algorithm

- Two-phase: Grouping of tasks, Assignment of task.
- Grouping of tasks using the level attribute. (lower level → higher priority) (same level → highest computation cost).
- Assign each task to the fastest available processor.
- Minimizes the sum of the task's computation cose and total communication costs with tasks in the previous levels.
- Disadvantages:
 - ▶ Time Complexity $O(v^2 \times q^2)$ with contention.



TASK-SCHEDULING ALGORITHM

Graph Attributes Used by HEFT and CPOP

Upward Rank

$$rank_u(n_{exit}) = \overline{w_{exit}}$$

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

- $rank_u(n_i)$: upward rank of n_i
- $\overline{w_{exit}}$: average execution cost of n_{exit}
- $succ(n_i)$: successor of n_i
- $\overline{c_{i,j}}$: average communication cost of edge (i, j)



TASK-SCHEDULING ALGORITHM

Graph Attributes Used by HEFT and CPOP

Downward Rank

$$rank_d(n_{entry}) = 0$$

$$rank_d(n_i) = \max_{n_j \in pred(n_i)} \{ rank_d(n_j) + \overline{w_j} + \overline{c_{i,j}} \}$$

- $rank_d(n_i)$: downward rank of n_i
- $pred(n_i)$: predecessor of n_i
- $\overline{w_j}$: average execution cost of n_j
- $\overline{c_{i,j}}$: average communication cost of edge (i, j)



HEFT

HETEROGENEOUS-EARLIEST-FINISH-TIME

Task Prioritizing Phase

- Priority of the tasks to be set to their upward rank value.
- Sort the tasks by decreasing order of upward ranks. Random Tie Breaking.
- Decreasing order provides topological ordering of the tasks.

Processor Selection Phase

- Earliest available time of processor is when it finishes the execution of its last assigned task.
- Insertion based policy.

Time complexity

- $O(eq)$ for e edges and q processors
- $O(v^2 \times q)$ for Dense Graph.

HEFT

HETEROGENEOUS-EARLIEST-FINISH-TIME

1. Set the computation costs of tasks and communication costs of edges with mean values.
2. Compute $rank_u$ for all tasks by traversing graph upward, starting from the exit task.
3. Sort the tasks in a scheduling list by nonincreasing order of $rank_u$ values.
4. **while** there are unscheduled tasks in the list **do**
5. Select the first task, n_i , from the list for scheduling.
6. **for** each processor p_k in the processor-set ($p_k \in Q$) **do**
7. Compute $EFT(n_i, p_k)$ value using the *insertion-based scheduling* policy.
8. Assign task n_i to the processor p_j that minimizes EFT of task n_i .
9. **endwhile**

Figure: HEFT algorithm.



Task Prioritizing Phase

- Priority of the tasks: summation of upward and downward ranks
- Entry task followed by immediate successor with highest priority is selected and marked as critical path task.
- Priority Queue is used to extract high priority task.
- Tie-breaking is done with first immediate successor with highest priority.



Processor Selection Phase

- Critical Path Processor minimizes the cumulative computation cost of the tasks on the critical path.
- Tasks on the critical path: scheduled on the critical path processor.
- other tasks: assigned to a processor which minimizes the earliest EFT.
- Insertion based policy.

Time complexity

- $O(eq)$ for e edges and q processors



CPOP

CRITICAL-PATH-ON-A-PROCESSOR

1. Set the computation costs of tasks and communication costs of edges with mean values.
2. Compute $rank_u$ of tasks by traversing graph upward, starting from the exit task.
3. Compute $rank_d$ of tasks by traversing graph downward, starting from the entry task.
4. Compute $priority(n_i) = rank_d(n_i) + rank_u(n_i)$ for each task n_i in the graph.
5. $|CP| = priority(n_{entry})$, where n_{entry} is the entry task.
6. $SET_{CP} = \{n_{entry}\}$, where SET_{CP} is the set of tasks on the critical path.
7. $n_k \leftarrow n_{entry}$.
8. **while** n_k is not the exit task **do**
9. Select n_j where $((n_j \in succ(n_k)) \text{ and } (priority(n_j) == |CP|))$.
10. $SET_{CP} = SET_{CP} \cup \{n_j\}$.
11. $n_k \leftarrow n_j$.
12. **endwhile**
13. Select the critical-path processor (p_{CP}) which minimizes $\sum_{n_i \in SET_{CP}} w_{i,j}, \forall p_j \in Q$.
14. Initialize the priority queue with the entry task.
15. **while** there is an unscheduled task in the priority queue **do**
16. Select the highest priority task n_i from priority queue.
17. **if** $n_i \in SET_{CP}$ **then**
18. Assign the task n_i on p_{CP} .
19. **else**
20. Assign the task n_i to the processor p_j which minimizes the $EFT(n_i, p_j)$.
21. Update the priority-queue with the successors of n_i , if they become ready tasks.
22. **endwhile**

Figure: CPOP algorithm.



PERFORMANCE EVALUATION METRICS

Comparison Metrics

- Schedule Length Ratio (SLR) (Makespan)

$$SLR = \frac{makespan}{\sum_{n_i \in CP_{MIN}} \min_{p_j \in Q} \{w_{i,j}\}}$$

- Speedup

$$speedup = \frac{\min_{p_j \in Q} \left\{ \sum_{n_i \in V} w_{i,j} \right\}}{makespan}$$

- Efficiency

$$efficiency = \frac{speedup}{q}$$

- Number of occurrences of better quality schedules
- Running time of the algorithms



PERFORMANCE EVALUATION

RANDOMLY GENERATED APPLICATION GRAPHS

Random Graph Generator: Input Parameters

- v : number of tasks in the graph
- α : shape parameter of graph
- out_degree : outdegree of a node
- CCR : communication to computation ratio
- β : range percentage of computation costs on processors



PERFORMANCE EVALUATION

RANDOMLY GENERATED APPLICATION GRAPHS

Range of Input Parameters

- SET_v : {20, 40, 60, 80, 100}
- SET_α : {0.5, 1.0, 2.0}
- SET_{out_degree} : {1, 2, 3, 4, 5, v }
- SET_{CCR} : {0.1, 0.5, 1.0, 5.0, 10.0}
- SET_β : {0.1, 0.25, 0.5, 0.75, 1.0}

$5 \times 3 \times 6 \times 5 \times 5 = 2250$ DAG types

$2250 \times 25 = 56250$ DAGs



PERFORMANCE RESULTS FOR RANDOM GRAPHS

- Average SLR performance ranking: HEFT, CPOP, DLS, MH, LMT
- Average speedup ranking: HEFT, DLS, CPOP=MH, LMT
- Running time ranking: HEFT, CPOP, MH, LMT, DLS
- Graph structure ranking:
 - ▶ $\alpha = 0.5$: HEFT, CPOP, MH, DLS, LMT
 - ▶ $\alpha = 1.0$: HEFT, CPOP=DLS, MH, LMT
 - ▶ $\alpha = 2.0$: HEFT, CPOP, DLS, MH, LMT
- CCR ranking:
 - ▶ $CCR \leq 1.0$: HEFT, CPOP, DLS, MH, LMT
 - ▶ $CCR > 1.0$: HEFT, CPOP, DLS, MH, LMT
- Number of best occurrences: HEFT, DLS, CPOP, MH, LMT



PERFORMANCE RESULTS

FOR REAL WORLD PROBLEMS APPLICATION GRAPHS

Gaussian Elimination Algorithm

- Matrix dimension = m (5 taken)
- Total number of tasks = $\frac{m^2+m+2}{2}$
- Structure of Graph is known.
- Avg SLR values: HEFT and DLS perform best ($q=5$, $m=5$ to 20).
- Efficiency: If number of processors > 8 (fixed m), then HEFT outperforms DLS.
- Running Time: DLS algorithm is the slowest.



PERFORMANCE RESULTS

FOR REAL WORLD PROBLEMS APPLICATION GRAPHS

Fast Fourier Transformation Algorithm

- Recursive, one-dimensional FFT Algorithm (four data points).
- Computation costs of same level tasks are equal.
- Communication costs of edges between two consecutive levels are same.
- Avg SLR values (various size of input points): HEFT outperforms others.
- Efficiency ($q=2$ to 32): HEFT and DLS gives most efficient schedules.
- Running Time (vary with number of data points and processors): DLS algorithm is highest cost.



PERFORMANCE RESULTS

FOR REAL WORLD PROBLEMS APPLICATION GRAPHS

Molecular Dynamics Code

- Number of tasks is fixed and structure of graph is known.
- CCR values ($q=6$):
- SLR ranking: **HEFT**, DLS, **CPOP**, MH, LMT
- DLS and LMT run three times slower than others.
- HEFT most practical and efficient.



ALTERNATE POLICIES FOR HEFT PHASES

TASK PRIORITIZING PHASE

priority(n_i) =

- A1: $rank_u(n_i) + rank_d(n_i)$
- A2: $rank_u(n_i) + \max_{n_j \in pred(n_i)} \{AFT(n_j)\}$
- A3: $rank_u(n_i) + \max_{n_j \in pred(n_i)} \{AFT(n_j) + c_{j,i}\}$

Observation

- Original priority policy gives better results than these alternates by 6 percent.



ALTERNATE POLICIES FOR HEFT PHASES

PROCESSOR SELECTION PHASE

critical_child(n_i) =

- B1: $\max_{n_c \in \text{succ}(n_i)} c_{i,c}$
- B2: $\max_{n_c \in \text{succ}(n_i)} \text{rank}_u(n_c)$
- B3: $\max_{n_c \in \text{succ}(n_i)} \{c_{i,c} + \text{rank}_u(n_c)\}$

Observation

- For small CCR graphs, original HEFT outperforms
- $3.0 \leq CCR < 6.0$: B1 outperforms original HEFT
- $CCR \geq 6.0$: B2 outperforms original HEFT and others alternates.



CONCLUSION

In Paper

- Two New Algorithm: HEFT and CPOP.
- HEFT outperforms the other algorithms.
- CPOP has given either better performance and better running time or comparable results with them.
- New Method for processor selection phase.

Planned Future

- Investigate the trade-off between the quality of schedules of the algorithms.
- Extend HEFT for rescheduling tasks in response to changes in processor and network loads.
- Extend for arbitrary-connected networks by considering link contention.

Thank You

