

Performance Models for Data Transfers: A Case Study with Molecular Chemistry Kernels

Suraj KUMAR¹, Lionel EYRAUD-DUBOIS², and
Sriram KRISHNAMOORTHY¹

¹Pacific Northwest National Laboratory, Richland, USA

²Inria Bordeaux – Sud-Ouest, France

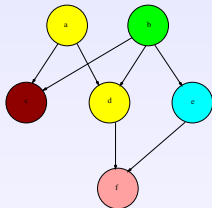


August 8, 2019



- Distributed memory systems are very common
- Rate of computation vs rate of data movement
- Focus: Avoid, Hide, and Minimize communication costs

Task Graphs and Runtime Systems

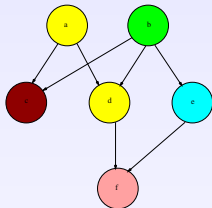


- Applications can be expressed as task graphs
- Vertices represent tasks
- Edges represent dependencies among tasks

Task Based Runtime Systems

- StarPU, QUARK, PaRSEC, StarSs, Legion, KAAPI
- May only see a set of ready (independent) tasks
- A memory node may require data from other memory nodes
- Order of data transfers such that communication-computation overlap is maximized

Task Graphs and Runtime Systems



- Applications can be expressed as task graphs
- Vertices represent tasks
- Edges represent dependencies among tasks

Task Based Runtime Systems

- StarPU, QUARK, PaRSEC, StarSs, Legion, KAAPI
- May only see a set of ready (independent) tasks
- A memory node may require data from other memory nodes
- Order of data transfers **between two memory nodes** such that communication-computation overlap is maximized

Problem Definition

- Communication and computation times are known in advance
 - ▶ Approximated based on number of computations and hardware details
 - ▶ Obtained from previous executions

Problem *DT*

- A set of tasks $ST = \{T_1, \dots, T_n\}$ is scheduled on a processing unit P with memory unit M of capacity C
- Input data for tasks of ST reside on another memory unit
- Tasks do not produce output data
- Tasks do not require intermediate memory
- A task uses an amount of memory in M from the start of its communication to the end of its computation

Given L , is there a feasible schedule S for ST such that makespan of S , $\mu(S) \leq L$?

Machine Flowshop Problem

- n machines and m tasks
- Each task contains exactly n operations
- i -th operation of a task must be processed on the i -th machine
- Each machine can perform at-max one operation at a time
- i -th operation of a task starts after the completion of $(i - 1)$ -th operation

The problem is to obtain the arrangement that achieves shortest possible makespan.

- Johnson provided an optimal algorithm for 2-machine flow shop problem
- Our problem DT adds one extra dimension (memory capacity) to 2-machine flowshop problem

1. Introduction

2. Problem Definition

- Unlimited Memory Capacity
- Limited Memory Capacity

3. Classes of Heuristics

- Static Order Based Strategies
- Dynamic Selection Based Strategies
- Static Order with Dynamic Corrections Based Strategies

4. Workloads & Hardware Configurations

5. Conclusion and Ongoing Work

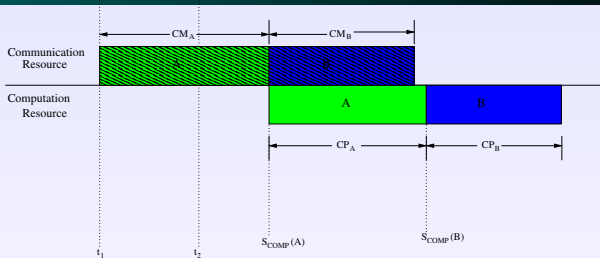
Unlimited Memory Capacity

Johnson's Algorithm

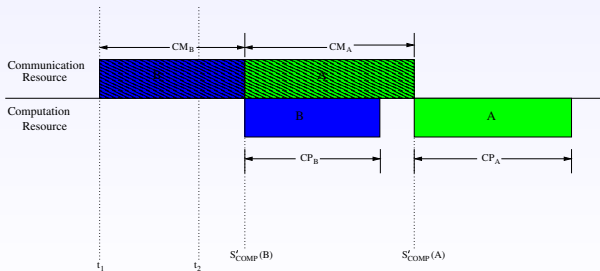
- 1: Divide ready tasks in two sets S_1 and S_2 . If computation time of a task T is not less than its communication time, then T is in S_1 otherwise in S_2 .
- 2: Sort S_1 in queue Q by non-decreasing communication times
- 3: Sort S_2 in queue Q' by non-increasing computation times
- 4: Append Q' to Q
- 5: $\tau_{\text{COMM}} \leftarrow 0$ {Available time of communication resource}
- 6: $\tau_{\text{COMP}} \leftarrow 0$ {Available time of computation resource}
- 7: **while** $Q \neq \emptyset$ **do**
- 8: Remove a task T from beginning of Q for processing
- 9: $S_{\text{COMM}}(T) \leftarrow \tau_{\text{COMM}}$
- 10: $S_{\text{COMP}}(T) \leftarrow \max(S_{\text{COMM}}(T) + \text{COMM}_T, \tau_{\text{COMP}})$
- 11: $\tau_{\text{COMM}} \leftarrow S_{\text{COMM}}(T) + \text{COMM}_T$
- 12: $\tau_{\text{COMP}} \leftarrow S_{\text{COMP}}(T) + \text{COMP}_T$
- 13: **end while**

- *OMIM* denotes the makespan of this algorithm

Approach for the optimality proof



Original Schedule



Swapped Schedule

1. Introduction

2. Problem Definition

- Unlimited Memory Capacity
- Limited Memory Capacity

3. Classes of Heuristics

- Static Order Based Strategies
- Dynamic Selection Based Strategies
- Static Order with Dynamic Corrections Based Strategies

4. Workloads & Hardware Configurations

5. Conclusion and Ongoing Work

- Memory is required only to store input data (from the definition of our problem)
- As Communication time \propto amount of communication, for each task:
 - ▶ Communication time = Amount of communication (for simplification)
 - ▶ Communication time = Amount of input data

Reduction Problem

Three Partition Problem (3PAR): Given a set of $3m$ integers $A = \{a_1, \dots, a_{3m}\}$, is there a partition of A into m triplets $TR_i = \{a_{i_1}, a_{i_2}, a_{i_3}\}$, such that $\forall i, a_{i_1} + a_{i_2} + a_{i_3} = b$, where $b = (1/m) \sum a_i$?

Definition of tasks in the reduction from 3PAR

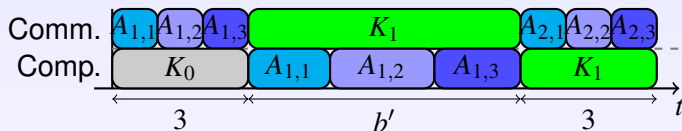
$$x = \max\{a_i : 1 \leq i \leq 3m\}$$

Task	Communication time	Computation time
K_0	0	3
K_1, \dots, K_{m-1}	$b' = b + 6x$	3
K_m	$b' = b + 6x$	0
$1 \leq i \leq 3m, A_i$	1	$a'_i = a_i + 2x$

Total memory capacity: $C = b' + 3$

Target makespan: $L = m(b' + 3)$

Pattern of Feasible Schedule

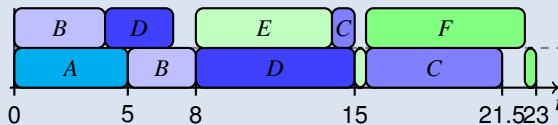


- Problem DT is NP-Complete.
- Our proof is inspired from work by Papadimitriou et al. (on 2-machine flowshop with limited buffer problem)

Order of Processing on Communication and Computation Resources in Optimal Schedules

Task	Memory Req	Comm Time	Comp Time
A	0	0	5
B	4	4	3
C	1	1	6
D	3	3	7
E	6	6	0.5
F	7	7	0.5

Common ordering on both resources



Different ordering on both resources

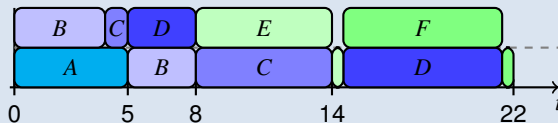


Table of Contents

1. Introduction

2. Problem Definition

- Unlimited Memory Capacity
- Limited Memory Capacity

3. Classes of Heuristics

- Static Order Based Strategies
- Dynamic Selection Based Strategies
- Static Order with Dynamic Corrections Based Strategies

4. Workloads & Hardware Configurations

5. Conclusion and Ongoing Work

- Static order based strategies
- Dynamic selection based strategies
- Static order with dynamic correction based strategies

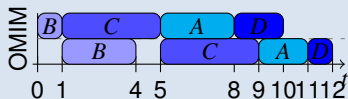
We consider common order on both resources for all of our heuristics.

- *order of optimal strategy infinite memory (OOSIM)*
- *increasing order of communication strategy (IOCMS)*
- *decreasing order of computation strategy (DOCPS)*
- *increasing order of communication plus computation strategy (IOCCS)*
- *decreasing order of communication plus computation strategy (DOCCS)*

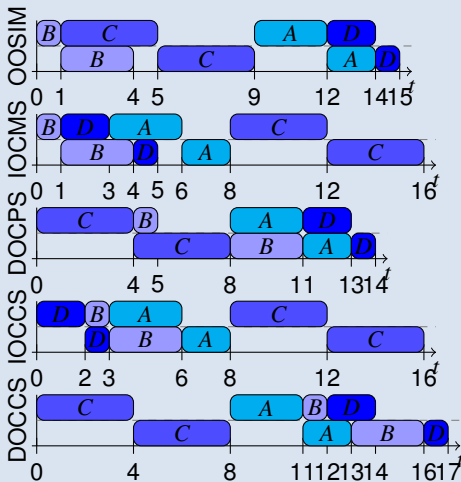
Static Order Based Strategies

Task	Memory Req	Comm Time	Comp Time
A	3	3	2
B	1	1	3
C	4	4	4
D	2	2	1

Unlimited Memory Capacity



Memory Capacity: 6

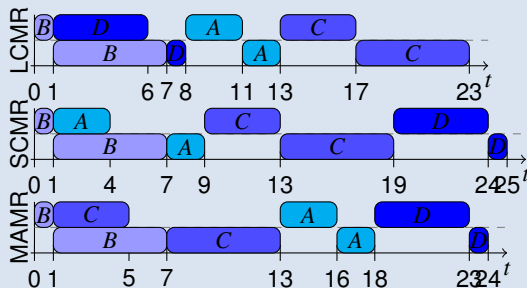


- *largest communication task respects memory restriction (LCMR)*
- *smallest communication task respects memory restriction (SCMR)*
- *maximum accelerated task respects memory restriction (MAMR)*

Dynamic Selection Based Strategies

Task	Memory Req	Comm Time	Comp Time
A	3	3	2
B	1	1	6
C	4	4	6
D	5	5	1

Memory Capacity = 6



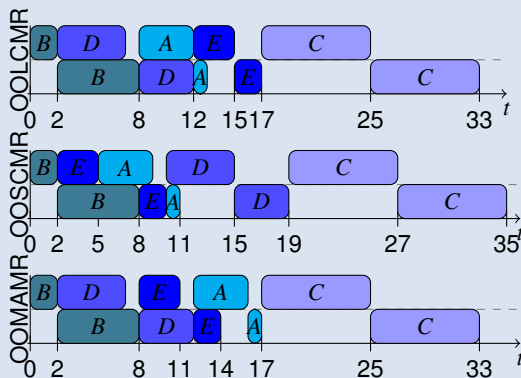
Static Order with Dynamic Corrections Based Strategies

- *optimal order infinite memory largest communication task respects memory restriction (OOLCMR)*
- *optimal order infinite memory smallest communication task respects memory restriction (OOSCMR)*
- *optimal order infinite memory maximum accelerated task respects memory restriction (OOMAMR)*

Static Order with Dynamic Corrections Based Strategies

Task	Memory Req	Comm Time	Comp Time
A	4	4	1
B	2	2	6
C	8	8	8
D	5	5	4
E	3	3	2

Memory Capacity=9



Favorable Situations for Heuristics

Heuristic	Favorable Situation
<i>OOSIM</i>	Memory capacity is not a restriction (Optimal)
<i>IOCMS</i>	Memory capacity is not a restriction and tasks are compute intensive (Optimal)
<i>DOCPS</i>	Memory capacity is not a restriction and tasks are communication intensive (Optimal)
<i>IOCCS</i>	Moderate memory capacity and most tasks are highly compute intensive
<i>DOCCS</i>	Moderate memory capacity and most tasks are highly communication intensive
<i>LCMR</i>	Limited memory capacity and significant percentage of tasks with large communication times are compute intensive
<i>SCMR</i>	Limited memory capacity and significant percentage of tasks with small communication times are compute intensive
<i>MAMR</i>	Limited memory capacity and significant percentage of all types of tasks
<i>OOLCMR</i>	Moderate memory capacity and significant percentage of slightly communication intensive tasks have large communication times
<i>OOSCMR</i>	Moderate memory capacity and significant percentage of slightly communication intensive tasks have small communication times
<i>OOMAMR</i>	Moderate memory capacity and significant percentage of all types of tasks

Gilmore-Gomory (GG)

- A classical algorithm for 2-machine no-wait flowshop problem
- Problem is represented by a graph
- An optimal sequence of vertices is obtained
- Does not take memory constraints into account

Bin Packing (BP)

- First-Fit algorithm for the bin packing problem
- Tasks added to the first bin in which they fit
- If no bin is found then a new bin is created
- Sequence made of all tasks from consecutive bins

Table of Contents

1. Introduction

2. Problem Definition

- Unlimited Memory Capacity
- Limited Memory Capacity

3. Classes of Heuristics

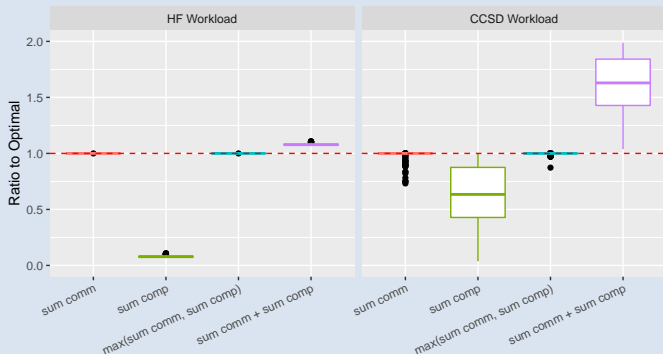
- Static Order Based Strategies
- Dynamic Selection Based Strategies
- Static Order with Dynamic Corrections Based Strategies

4. Workloads & Hardware Configurations

5. Conclusion and Ongoing Work

Molecular Chemistry Kernels

- Hartree–Fock (HF) with SiOSi molecules and Coupled Cluster Singles Doubles (CCSD) with Uracil molecules
- Tensor operations: transpose and contraction



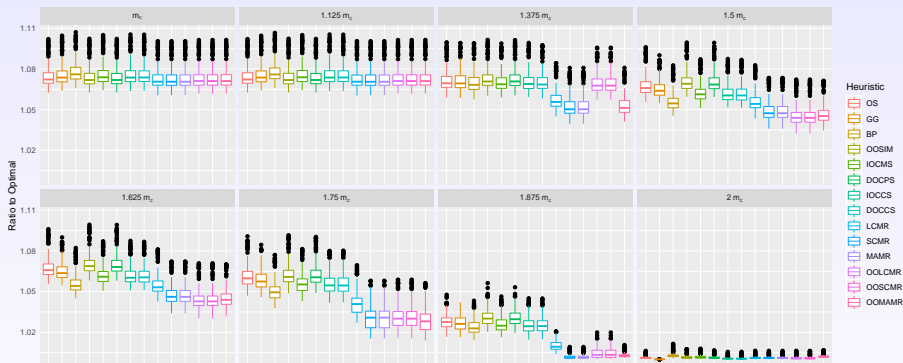
Configuration Parameters to Obtain Traces

- 10 nodes of Cascade machine
- Each node contains 16 Intel Xeon E5-2670 cores
- Double precision version of HF and CCSD of NWChem
- NWChem takes advantages of Global Arrays (GA)
- 150 processes for each application, 300-800 tasks for each process

Simulation Parameters

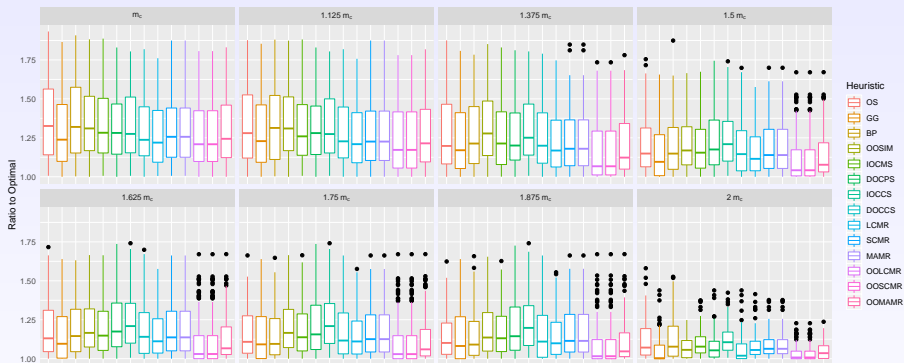
- m_c : minimum memory capacity requirement to execute all tasks of an application
- Evaluation criteria for a heuristic H , $r(H) = \frac{\text{makespan of } H}{OMIM}$ (lower values are better)
- All data transfers between the local memory of each process and the GA memory take the same route

HF Performance with $m_c = 176KB$.



- Dynamic strategies are best suited for limited memory capacities
- Static order with dynamic correction variants outperform others for moderate memory capacities

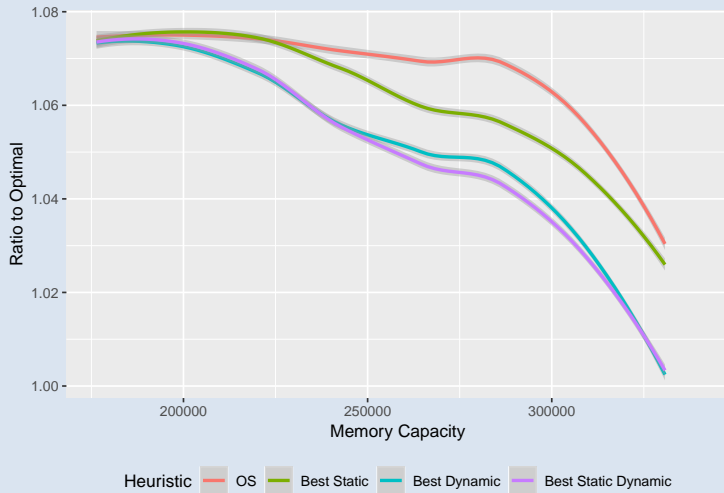
CCSD performance with $m_c = 1.8GB$



- Highly heterogeneous tasks
- Static order with dynamic correction variants outperform others in most cases

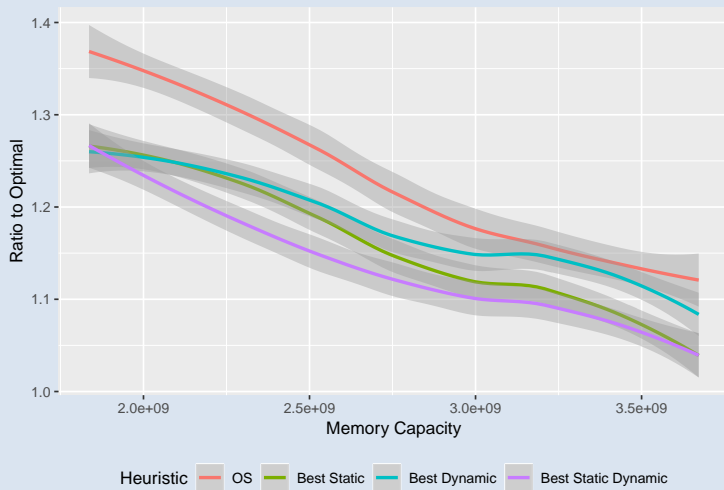
Best variants of all categories

HF Performance



Best variants of all categories

CCSD Performance



Implementation challenges

- Impact of Congestion on communication/computation times
- Output data of a task
- Routes between two memory nodes
- Communications from multiple memory nodes can happen at the same time

Table of Contents

1. Introduction

2. Problem Definition

- Unlimited Memory Capacity
- Limited Memory Capacity

3. Classes of Heuristics

- Static Order Based Strategies
- Dynamic Selection Based Strategies
- Static Order with Dynamic Corrections Based Strategies

4. Workloads & Hardware Configurations

5. Conclusion and Ongoing Work

Conclusion:

- Problem of determining the optimal order of data transfers is NP-complete
- Our heuristics achieve significant overlap for HF and CCSD applications

Ongoing work:

- Evaluation of our heuristics in the context of accelerators
- Implementation of the proposed heuristics
- Automatic selection of the best heuristic
- Model bandwidth sharing to support multiple simultaneous communications