Efficient scheduling of DAGs under bounded memory

Loris Marchal Joint work with Bertrand Simon & Frédéric Vivien

CNRS and University of Lyon

Journée Mao — Besançon 2019





Modeling scientific applications as task graphs

- Scientific applications divided into rather independent modules (tasks)
- Tasks linked through data dependencies
- Directed Acyclic Graph of tasks



Multifrontal sparse matrix factorization over runtimes





- Task graph: tree (with dependencies towards the root)
- Large temporary data
- Memory becomes a bottleneck
- Schedule trees with limited memory

Consider a simple task graph



▲ロト ▲圖ト ▲ヨト ▲ヨト ニヨー のへで

- Consider a simple task graph
- Tasks have durations and memory demands



- Consider a simple task graph
- Tasks have durations and memory demands



time

- Consider a simple task graph
- Tasks have durations and memory demands



▲ロ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ● ○ ○ ○

Peak memory: maximum memory usage

- Consider a simple task graph
- Tasks have durations and memory demands



- Peak memory: maximum memory usage
- Trade-off between peak memory and makespan

Going back to sequential processing

- Temporary data require memory
- Scheduling influences the peak memory



イロト イポト イヨト イヨト

Going back to sequential processing

- Temporary data require memory
- Scheduling influences the peak memory



イロト 不得下 不良下 不良下

э

Going back to sequential processing

- Temporary data require memory
- Scheduling influences the peak memory



Known results on sequential memory-aware scheduling:

- Optimal algorithm for trees [Liu, 1897], SP-graphs [Kayaaslan et al., 2018]
- General graphs with unit size weights: pebble game [Sethi and Ullman, 1973], PSPACE complete [Gilbert et al., 1980]

Today's focus

Dynamic scheduling of general graphs:



On shared-memory platforms:



Objective: limit memory consumption while allowing parallelism

Outline

Model and maximum parallel memory

Memory model Maximum parallel memory/maximal topological cut

Efficient scheduling with bounded memory

Problem definition Complexity Heuristics Simulation results

Conclusion

Outline

Model and maximum parallel memory

Memory model Maximum parallel memory/maximal topological cut

Efficient scheduling with bounded memory

Problem definition Complexity Heuristics Simulation results

Conclusion

Task graphs with:

▶ Vertex weights (*w_i*): task (estimated) durations

▶ Edge weights (*m_{i,j}*): data sizes

Task graphs with:

- ▶ Vertex weights (*w_i*): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory



Task graphs with:

- ▶ Vertex weights (*w_i*): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory



Task graphs with:

- ▶ Vertex weights (*w_i*): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory



Task graphs with:

- ▶ Vertex weights (*w_i*): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory



Task graphs with:

- ▶ Vertex weights (*w_i*): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory



Task graphs with:

- ▶ Vertex weights (*w_i*): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory

Emulation of other memory behaviours:

Inputs + outputs allocated during task: duplicate nodes



Task graphs with:

- Vertex weights (w_i): task (estimated) durations
- ▶ Edge weights (*m_{i,j}*): data sizes

Simple memory model: at the beginning of a task

- Inputs are freed (instantaneously)
- Outputs are allocated

At the end of a task: outputs stay in memory

Emulation of other memory behaviours:

- Inputs + outputs allocated during task: duplicate nodes
- ► Shared data: output data of A used for both B and C:



Computing the maximum memory peak

Topological cut: (S, T) with:

- S include the source node, T include the target node
- No edge from T to S
- Weight of the cut = weight of all edges from S to T



Any topological cut corresponds to a possible state when all node in S are completed or being processed.

Two equivalent questions (in our model):

What is the maximum memory of any parallel execution?

What is the topological cut with maximum weight?

Computing the maximum topological cut

Literature:

- Lots of studies of various cuts in non-directed graphs ([Diaz,2000] on Graph Layout Problems)
- Minimum cut is polynomial on both directed/non-directed graphs
- Maximum cut NP-complete on both directed/non-directed graphs ([Karp 1972] for non-directed, [Lampis 2011] for directed ones)
- Not much for topological cuts

Theorem.

Computing the maximum topological cut of a DAG can be done in polynomial time.

 Consider one classical LP formulation for finding a minimum cut:

$$egin{aligned} \min \sum_{(i,j)\in E} m_{i,j}d_{i,j} \ orall (i,j)\in E, \quad d_{i,j}\geq p_i-p_j \ orall (i,j)\in E, \quad d_{i,j}\geq 0 \ p_s=1, \quad p_t=0 \end{aligned}$$

 Consider one classical LP formulation for finding a minimum cut:

$$\min \sum_{(i,j)\in E} m_{i,j}d_{i,j}$$
 $orall (i,j)\in E, \quad d_{i,j}\geq p_i-p_j$
 $orall (i,j)\in E, \quad d_{i,j}\geq 0$
 $p_s=1, \quad p_t=0$

▲□▶ ▲□▶ ▲ 臣▶ ★ 臣▶ 三臣 … のへぐ

► Integer solution ⇔ topological cut

 Consider one classical LP formulation for finding a minimum cut:

$$\max \sum_{(i,j)\in E} m_{i,j}d_{i,j}$$
$$\forall (i,j)\in E, \quad d_{i,j}=p_i-p_j$$
$$\forall (i,j)\in E, \quad d_{i,j}\geq 0$$
$$p_s=1, \quad p_t=0$$

- ► Integer solution ⇔ topological cut
- Then change the optimization direction (min \rightarrow max)

 Consider one classical LP formulation for finding a minimum cut:

$$\max \sum_{(i,j)\in E} m_{i,j}d_{i,j}$$
$$\forall (i,j)\in E, \quad d_{i,j}=p_i-p_j$$
$$\forall (i,j)\in E, \quad d_{i,j}\geq 0$$
$$p_s=1, \quad p_t=0$$

- ► Integer solution ⇔ topological cut
- Then change the optimization direction (min \rightarrow max)

► Draw w uniformly in]0,1[, define the cut such that $S_w = \{i \mid p_i > w\}, \quad T_w = \{i \mid p_i \le w\}$

 Consider one classical LP formulation for finding a minimum cut:

$$\max \sum_{(i,j)\in E} m_{i,j}d_{i,j}$$
$$\forall (i,j)\in E, \quad d_{i,j}=p_i-p_j$$
$$\forall (i,j)\in E, \quad d_{i,j}\geq 0$$
$$p_s=1, \quad p_t=0$$

- ► Integer solution ⇔ topological cut
- Then change the optimization direction (min \rightarrow max)
- ▶ Draw *w* uniformly in]0,1[, define the cut such that $S_w = \{i \mid p_i > w\}, \quad T_w = \{i \mid p_i \le w\}$
- ► Expected cost of this cut = M^{*} (opt. rational solution)

 Consider one classical LP formulation for finding a minimum cut:

$$\max \sum_{(i,j)\in E} m_{i,j}d_{i,j}$$
$$\forall (i,j)\in E, \quad d_{i,j}=p_i-p_j$$
$$\forall (i,j)\in E, \quad d_{i,j}\geq 0$$
$$p_s=1, \quad p_t=0$$

- Integer solution \Leftrightarrow topological cut
- Then change the optimization direction (min \rightarrow max)
- ► Draw *w* uniformly in]0,1[, define the cut such that $S_w = \{i \mid p_i > w\}, \quad T_w = \{i \mid p_i \le w\}$
- ► Expected cost of this cut = M^{*} (opt. rational solution)
- ► All cuts with random w have the same cost $M_{\mathbb{P}}^*$, \mathbb{P} , \mathbb{P} , \mathbb{P}

- Dual problem: Min-Flow (larger than all edge weights)
- Idea: use an optimal algorithm for Max-Flow

Algorithm sketch

- 1. Build a large flow F on the graph G
- 2. Consider G^{diff} with edge weights $F_{i,j} m_{i,j}$
- 3. Compute a maximum flow maxdiff in G^{diff}
- 4. F maxdiff is a minimum flow in G
- 5. Residual graph \rightarrow maximum topological cut



Complexity: same as maximum flow, e.g., $O(|V|^2|E|)$

- Dual problem: Min-Flow (larger than all edge weights)
- Idea: use an optimal algorithm for Max-Flow

Algorithm sketch1. Build a large flow F on the graph G2. Consider G^{diff} with edge weights $F_{i,j} - m_{i,j}$ 3. Compute a maximum flow maxdiff in
 G^{diff} 4. F - maxdiff is a minimum flow in G5. Residual graph \rightarrow maximum

Complexity: same as maximum flow, e.g., $O(|V|^2|E|)$

topological cut

- Dual problem: Min-Flow (larger than all edge weights)
- Idea: use an optimal algorithm for Max-Flow

Algorithm sketch $F_{i,j}$ 1. Build a large flow F on the graph G2. Consider G^{diff} with edge weights $F_{i,j} - m_{i,j}$ 3. Compute a maximum flow maxdiff in
 G^{diff} 4. F - maxdiff is a minimum flow in G5. Residual graph \rightarrow maximum

Complexity: same as maximum flow, e.g., $O(|V|^2|E|)$

topological cut

- Dual problem: Min-Flow (larger than all edge weights)
- Idea: use an optimal algorithm for Max-Flow

Algorithm sketch

- 1. Build a large flow F on the graph G
- 2. Consider G^{diff} with edge weights $F_{i,j} m_{i,j}$
- 3. Compute a maximum flow *maxdiff* in G^{diff}
- 4. F maxdiff is a minimum flow in G
- 5. Residual graph \rightarrow maximum topological cut





- Dual problem: Min-Flow (larger than all edge weights)
- Idea: use an optimal algorithm for Max-Flow

Algorithm sketch

- 1. Build a large flow F on the graph G
- 2. Consider G^{diff} with edge weights $F_{i,j} m_{i,j}$
- 3. Compute a maximum flow *maxdiff* in G^{diff}
- 4. F maxdiff is a minimum flow in G
- 5. Residual graph \rightarrow maximum topological cut



Complexity: same as maximum flow, e.g., $O(|V|^2|E|)$

Outline

Model and maximum parallel memory

Memory model Maximum parallel memory/maximal topological cut

Efficient scheduling with bounded memory

Problem definition Complexity Heuristics Simulation results

Conclusion

Coping with limiting memory

Problem:

- Limited available memory M
- Allow use of dynamic schedulers
- Avoid running out of memory
- Keep high level of parallelism (as much as possible)

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

Coping with limiting memory

Problem:

- Limited available memory M
- Allow use of dynamic schedulers
- Avoid running out of memory
- Keep high level of parallelism (as much as possible)

Our solution:

 Add edges to guarantee that any parallel execution stays below M fictitious dependencies to reduce maximum memory

 $M = \pm 0$ (F) (F)

Minimize the obtained critical path



Coping with limiting memory

Problem:

- Limited available memory M
- Allow use of dynamic schedulers
- Avoid running out of memory
- Keep high level of parallelism (as much as possible)

Our solution:

 Add edges to guarantee that any parallel execution stays below M fictitious dependencies to reduce maximum memory

 $M = 40 \quad \text{abs} \quad a \to b = b \quad a \to b \to b$

Minimize the obtained critical path



Definition (PartialSerialization).

Given a DAG G = (V, E) and a bound M, find a set of new edges E' such that $G' = (V, E \cup E')$ is a DAG, $MaxMem(G') \le M$ and CritPath(G') is minimized.

Theorem.

PartialSerialization is NP-hard in the stronge sense.

NB: stays NP-hard if we are given a sequential schedule σ of G which uses at most a memory M.

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ つ ・

Heuristic solutions for PARTIALSERIALIZATION

Framework:

(inspired by [Sbîrlea et al. 2014])

- 1. Compute a max. top. cut (S, T)
- 2. If weight $\leq M$: succeeds
- Add edge (u, v) with u ∈ T, v ∈ S without creating cycles; or fail



・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ つ ・

4. Goto Step 1

Several heuristic choices for Step 3:

MinLevels does not create a large critical path

RespectOrder follows a precomputed memory-efficient schedule, always succeeds

MaxSize targets nodes dealing with large data MaxMinSize variant of MaxSize

Simulations: dense random graphs (25, 50, 100 nodes)



- ► x: memory (0 = DFS, 1 = MaxTopCut) median ratio MaxTopCut / DFS ≈ 1.3
- y: CP / original CP \rightarrow lower is better
- MinLevels performs best

Simulations: sparse random graphs (25, 50, 100 nodes)



- ► x: memory (0 = DFS, 1 = MaxTopCut) median ratio MaxTopCut / DFS ≈ 2
- y: CP / original CP \rightarrow lower is better
- MinLevels performs best, but might fail

Simulations – Pegasus workflows (LIGO 100 nodes)



- ▶ Median ratio MaxTopCut / DFS ≈ 20
- MinLevels performs best, RespectOrder always succeeds

Memory divided by 5 for CP multiplied by 3

Simulations – Pegasus workflows (LIGO 100 nodes)



- ▶ Median ratio MaxTopCut / DFS ≈ 20
- MinLevels performs best, RespectOrder always succeeds

Memory divided by 5 for CP multiplied by 3

Outline

Model and maximum parallel memory

Memory model Maximum parallel memory/maximal topological cut

Efficient scheduling with bounded memory

Problem definition Complexity Heuristics Simulation results

Conclusion

Conclusion

Contributions:

- Maximum (parallel) memory of DAG:
 Weight of the maximum topological cut
- Bound memory usage by adding of fictitious dependencies

Take-away messages:

- Trade-off between parallelism and maximum memory on DAGs may be studied through cuts/flows.
- Apply to any DAG scheduling problem where storage matters
- Renewal in task graph scheduling: task-based runtime systems

- Such as: ParSec, StarPU, XKaapi
- Programmers write task code + dependencies
- Mapping/scheduling at runtime
- Need for lightweight online schedulers