

CUDA accelerated LTL Model Checking

Jiří Barnat, Luboš Brim, **Milan Češka**, and Tomáš Lamr

Faculty of Informatics, Masaryk University
Brno, Czech Republic

based on [Barnat et al. ICPADS'09] and [Barnat et al. PDMC'09]

Mathematical and Engineering Methods in Computer Science
(MEMICS'09)

13.11, 2009

① Introduction

Motivation and contribution

LTL Model Checking

Maximal Accepting Predecessor algorithm

② CUDA Accelerated MAP Algorithm

Reformulation of MAP algorithm

Compact CSR representation

③ Experimental evaluation

Impact of graph reduction

The overall run-times and speedup

④ Conclusion and Future Work

Motivation

Formal verification

- critical system - any mistake may have fatal consequences
- testing is insufficient - can not guarantee correctness
- formal methods can prove or disprove correctness of the system

Model Checking

- fully automated approach to the formal verification
- state space explosion problem
- possible solution:
 - symbolic representation
 - reduction techniques
 - **platform-dependent verification**

Motivation

Formal verification

- critical system - any mistake may have fatal consequences
- testing is insufficient - can not guarantee correctness
- formal methods can prove or disprove correctness of the system

Model Checking

- fully automated approach to the formal verification
- state space explosion problem
- possible solution:
 - symbolic representation
 - reduction techniques
 - **platform-dependent verification**

Platform-dependent verification

DiVinE

- Explicit Parallel LTL Model Checker
- Focuses on full utilization of available HW power

Successful stories on the following parallel platforms

- Clusters, Grids
DiVinE Cluster
- Multi-core workstations
DiVinE Multi-Core

Many-core architectures

- Shared-memory setting, many computing cores
- Parallel computing platform of the future?
- Widely accessible due to GP GPU devices
- Not yet covered by DiVinE

Platform-dependent verification

DiVinE

- Explicit Parallel LTL Model Checker
- Focuses on full utilization of available HW power

Successful stories on the following parallel platforms

- Clusters, Grids
 DiVinE Cluster
- Multi-core workstations
 DiVinE Multi-Core

Many-core architectures

- Shared-memory setting, many computing cores
- Parallel computing platform of the future?
- Widely accessible due to GP GPU devices
- Not yet covered by DiVinE

Platform-dependent verification

DiVinE

- Explicit Parallel LTL Model Checker
- Focuses on full utilization of available HW power

Successful stories on the following parallel platforms

- Clusters, Grids
 DiVinE Cluster
- Multi-core workstations
 DiVinE Multi-Core

Many-core architectures

- Shared-memory setting, many computing cores
- Parallel computing platform of the future?
- Widely accessible due to GP GPU devices
- Not yet covered by DiVinE

Platform-dependent verification

DiVinE

- Explicit Parallel LTL Model Checker
- Focuses on full utilization of available HW power

Successful stories on the following parallel platforms

- Clusters, Grids
 DiVinE Cluster
- Multi-core workstations
 DiVinE Multi-Core

Many-core architectures

- Shared-memory setting, many computing cores
- Parallel computing platform of the future?
- Widely accessible due to GP GPU devices
- Not yet covered by DiVinE

Platform-dependent verification

DiVinE

- Explicit Parallel LTL Model Checker
- Focuses on full utilization of available HW power

Successful stories on the following parallel platforms

- Clusters, Grids
DiVinE Cluster
- Multi-core workstations
DiVinE Multi-Core

Many-core architectures

- Shared-memory setting, many computing cores
- Parallel computing platform of the future?
- Widely accessible due to GP GPU devices
- Not yet covered by DiVinE

Platform-dependent verification

DiVinE

- Explicit Parallel LTL Model Checker
- Focuses on full utilization of available HW power

Successful stories on the following parallel platforms

- Clusters, Grids
 DiVinE Cluster
- Multi-core workstations
 DiVinE Multi-Core

Many-core architectures

- Shared-memory setting, many computing cores
- Parallel computing platform of the future?
- Widely accessible due to GP GPU devices
- Not yet covered by DiVinE

Many-core architecture

NVIDIA CUDA Technology

- Many-core architecture
- Hundreds of computing cores
- HW support for thousands of computing threads
- **Requires quite specific memory-usage pattern**
- Incompatible with random memory access (hashing)

Need for many-core algorithms

- Straightforward adaptation of multi-core algorithms to many-core architecture is inefficient
- Also the case of parallel LTL Model Checking

Many-core architecture

NVIDIA CUDA Technology

- Many-core architecture
- Hundreds of computing cores
- HW support for thousands of computing threads
- **Requires quite specific memory-usage pattern**
- Incompatible with random memory access (hashing)

Need for many-core algorithms

- Straightforward adaptation of multi-core algorithms to many-core architecture is inefficient
- Also the case of parallel LTL Model Checking

Contribution

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration [Barnat et al. ICPADS'09].

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking [Barnat et al. PDMC'09].

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Sketch of possible solutions

Contribution

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration [Barnat et al. ICPADS'09].

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking [Barnat et al. PDMC'09].

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Sketch of possible solutions

Contribution

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration [Barnat et al. ICPADS'09].

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking [Barnat et al. PDMC'09].

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Sketch of possible solutions

Contribution

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration [Barnat et al. ICPADS'09].

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking [Barnat et al. PDMC'09].

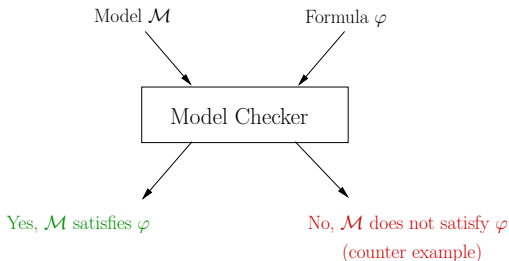
Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Sketch of possible solutions

LTL Model Checking

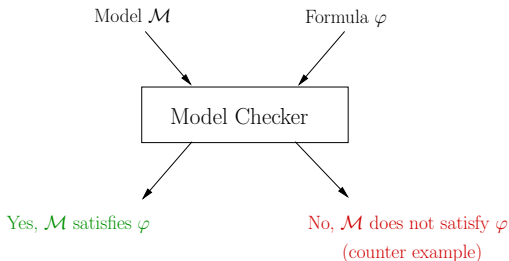
- inspected system \rightarrow suitable model M
- required property of the system \rightarrow formula φ in Linear Temporal Logic (LTL)
- problem : **Does the model M satisfy the formula φ ?**



- solution: automata based approach
 - reduction on accepting cycle detection

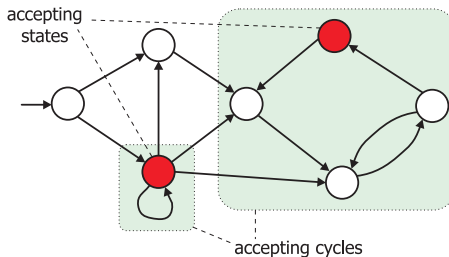
LTL Model Checking

- inspected system \rightarrow suitable model M
- required property of the system \rightarrow formula φ in Linear Temporal Logic (LTL)
- problem : **Does the model M satisfy the formula φ ?**



- solution: automata based approach
 - reduction on **accepting cycle detection**

Algorithms for accepting cycle detection



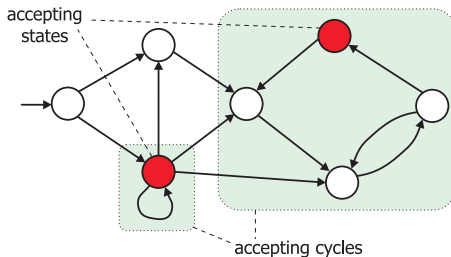
Optimal time complexity but hard to parallelize:

- Nested DFS
- Tarjan's SCC decomposition

Unoptimal time complexity but easy to parallelize:

- OWCTY
- **Maximal accepting predecessor (MAP)**

Algorithms for accepting cycle detection



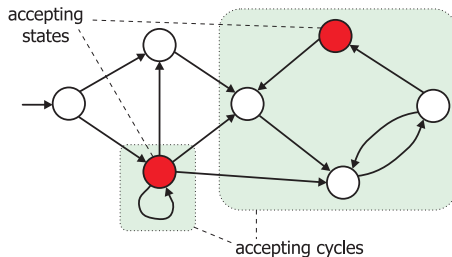
Optimal time complexity but hard to parallelize:

- Nested DFS
- Tarjan's SCC decomposition

Unoptimal time complexity but easy to parallelize:

- OWCTY
- **Maximal accepting predecessor (MAP)**

Algorithms for accepting cycle detection



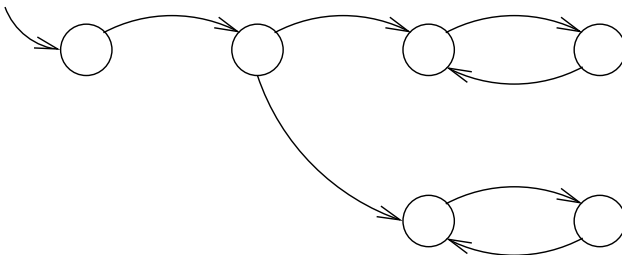
Optimal time complexity but hard to parallelize:

- Nested DFS
- Tarjan's SCC decomposition

Unoptimal time complexity but easy to parallelize:

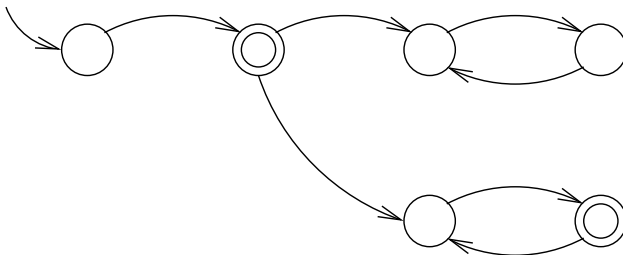
- OWCTY
- **Maximal accepting predecessor (MAP)**

Algorithm MAP



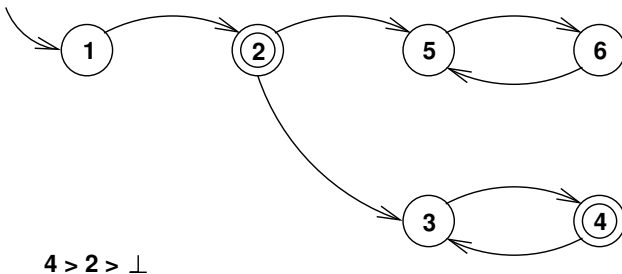
Graph corresponding to the state space.

Algorithm MAP



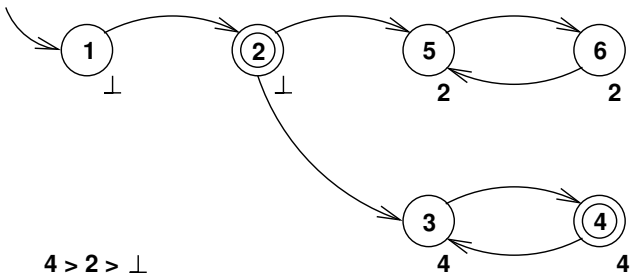
Accepting vertices, accepting cycle.

Algorithm MAP



Vertex ordering.

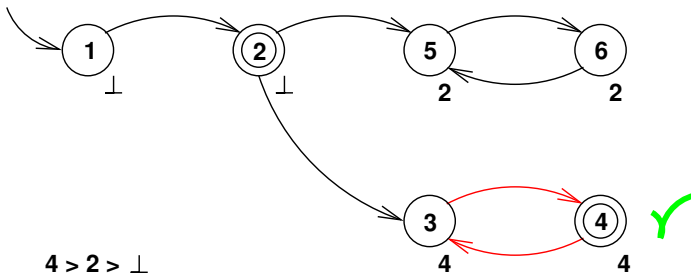
Algorithm MAP



Maximal Accepting Predecessor (MAP)

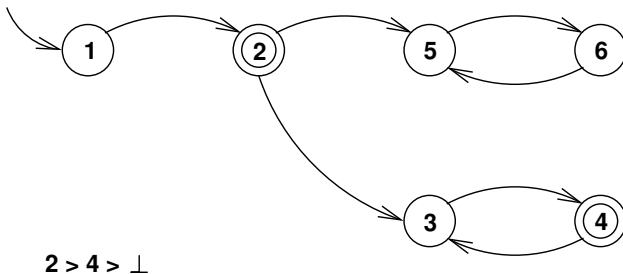
$$\text{map}(v) = \max\{\perp, u \mid (u, v) \in E^+ \wedge \mathcal{A}(u)\}$$

Algorithm MAP



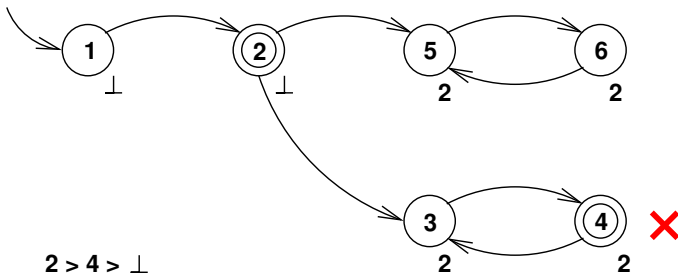
$map(v) = v \implies$ accepting cycle

Algorithm MAP



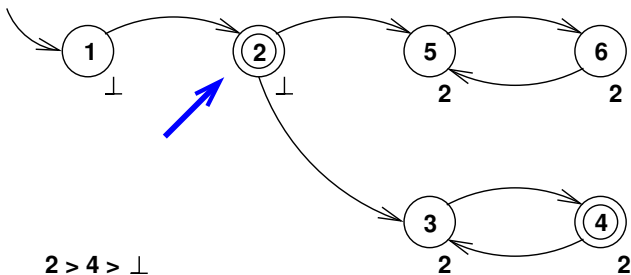
What if $2 > 4$?

Algorithm MAP



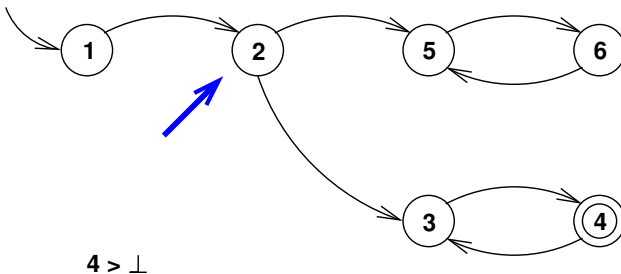
Accepting cycle undetected.

Algorithm MAP



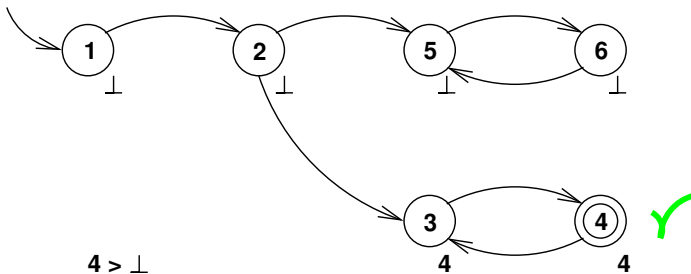
If no accepting cycle is found, then maximal accepting vertices
 cannot be part of an accepting cycle.

Algorithm MAP



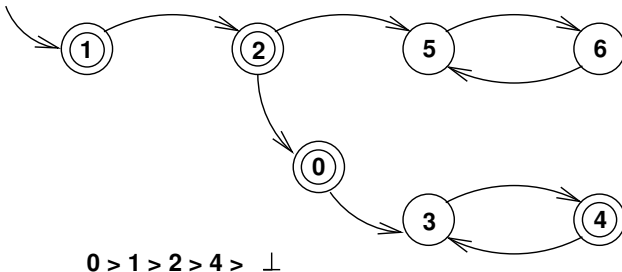
Maximal accepting vertices marked as non-accepting.

Algorithm MAP

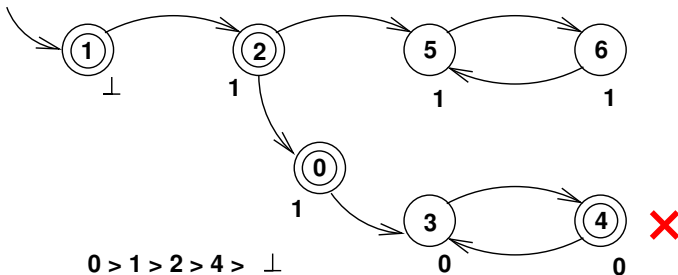


Repeat until accepting cycle is found or there are no accepting vertices.

Algorithm MAP - partitioning into subgraphs

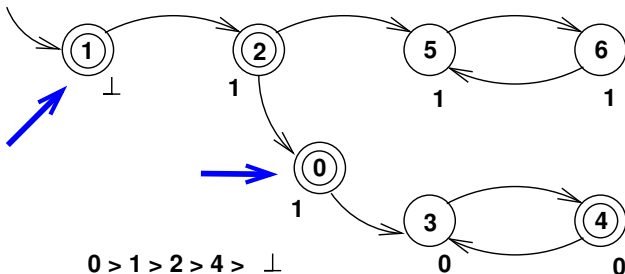


Algorithm MAP - partitioning into subgraphs



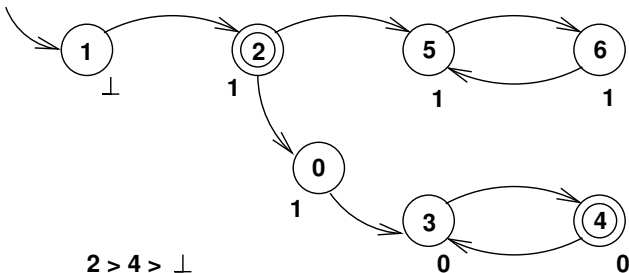
Accepting cycle undetected.

Algorithm MAP - partitioning into subgraphs



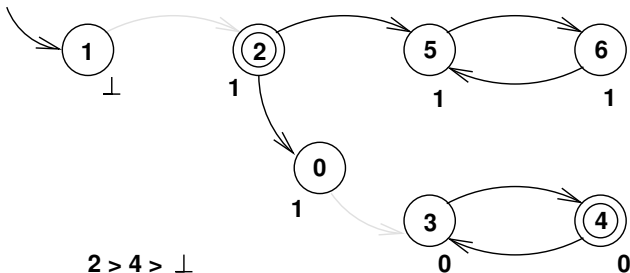
If no accepting cycle is found, then maximal accepting vertices
 cannot be part of an accepting cycle.

Algorithm MAP - partitioning into subgraphs



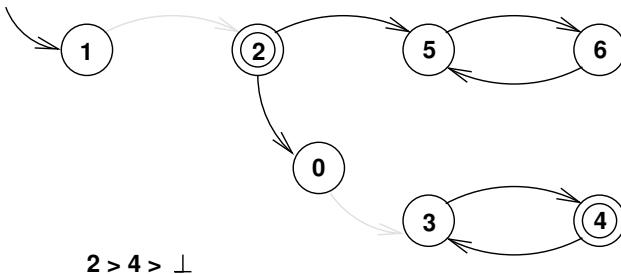
Maximal accepting vertices marked as non-accepting.

Algorithm MAP - partitioning into subgraphs



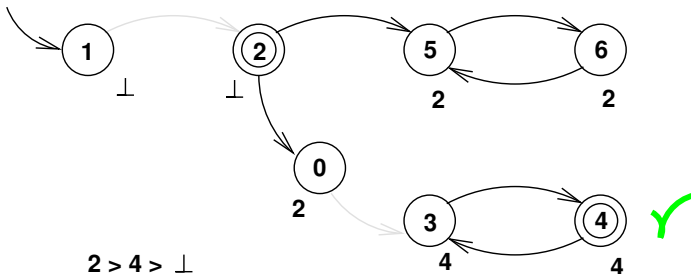
Partitioning into subgraphs according to the *map* values.

Algorithm MAP - partitioning into subgraphs



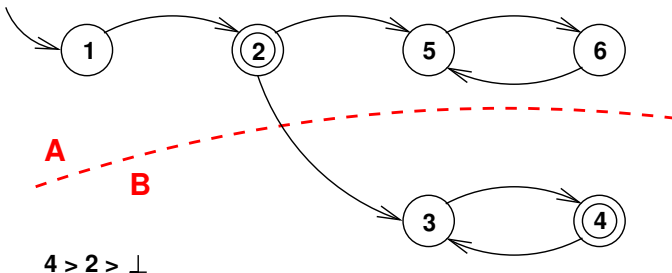
No cycle of the original graph maps to multiple partitions.

Algorithm MAP - partitioning into subgraphs



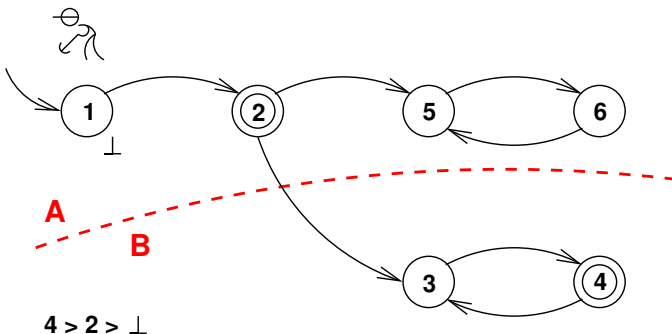
Earlier detection of accepting cycle.

Computing values of map – multi-cores



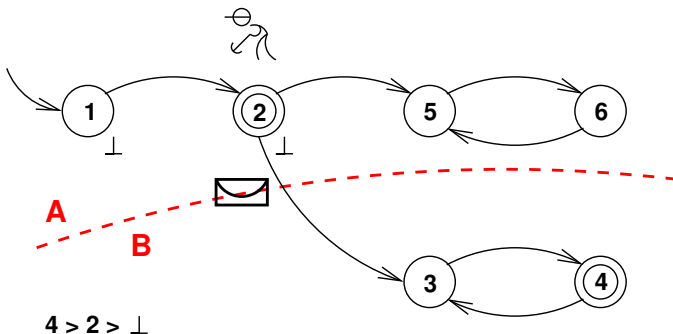
Graph partitioning.

Computing values of map – multi-cores



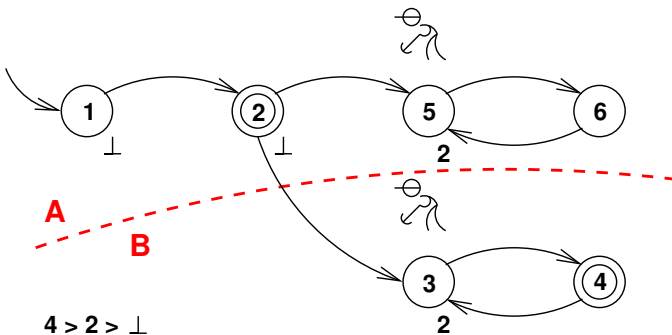
Each core processes own vertices.

Computing values of map – multi-cores



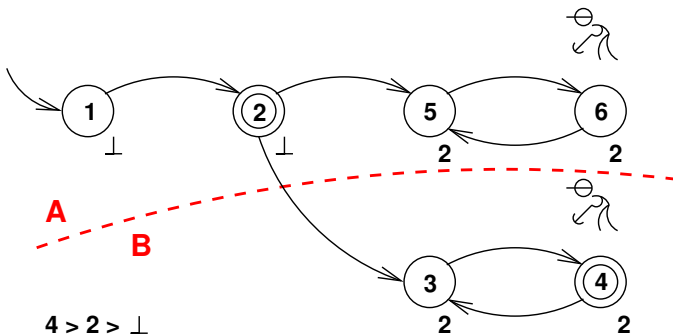
Non local vertices are sent to the owners.

Computing values of map – multi-cores



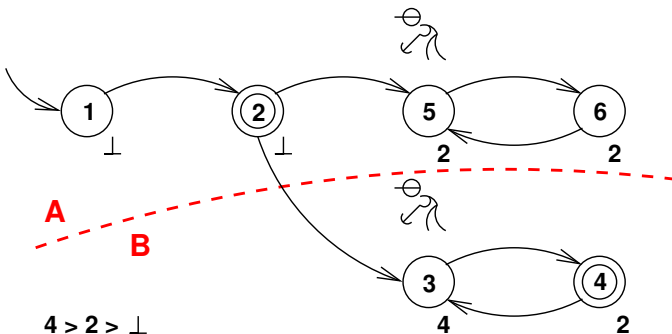
Parallel processing of vertices.

Computing values of map – multi-cores



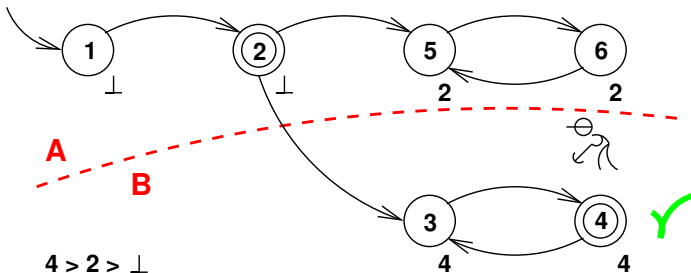
Parallel processing of vertices.

Computing values of map – multi-cores



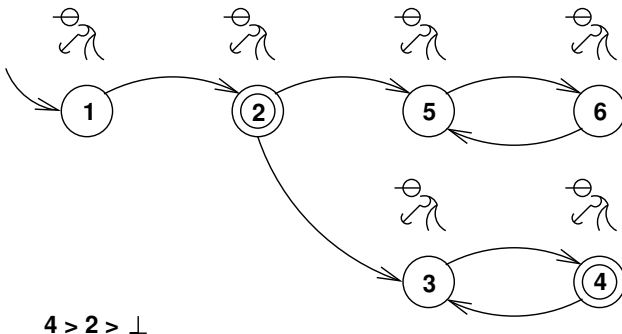
Parallel processing of vertices.

Computing values of map – multi-cores



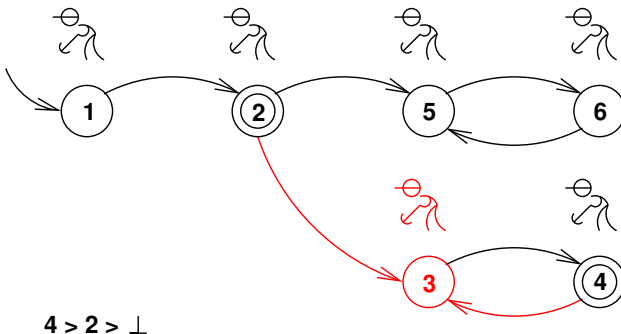
Accepting cycle found.

Computing values of map – many-cores



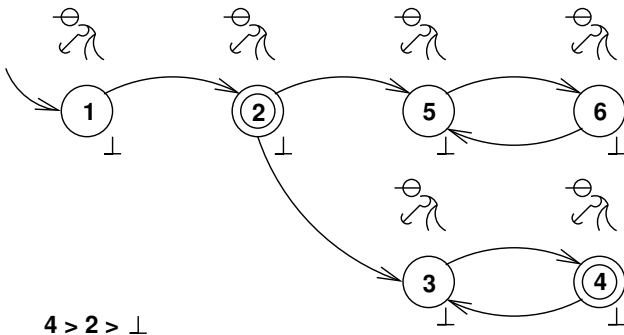
One thread per vertex.

Computing values of map – many-cores



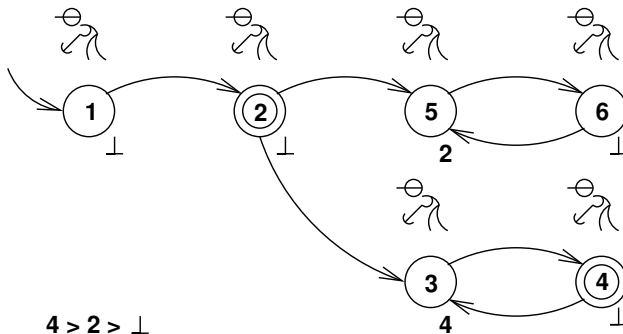
Each thread processes all incoming edges.

Computing values of map – many-cores



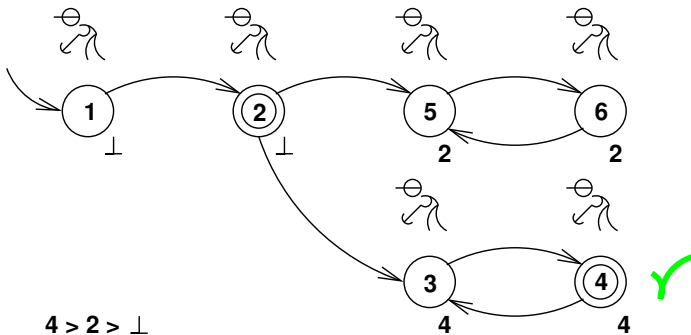
Threads proceed simultaneously.

Computing values of map – many-cores



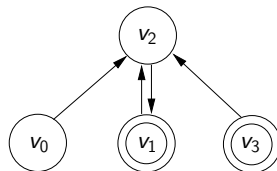
Threads proceed simultaneously.

Computing values of map – many-cores



Accepting cycle found.

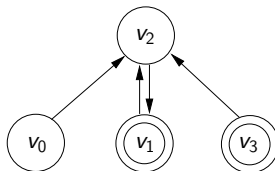
MAP Algorithm as Matrix-Vector Product



$$\begin{array}{c}
 v_0 \ v_1 \ v_2 \ v_3 \\
 v_0 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\
 v_1 \\
 v_2 \\
 v_3
 \end{array}$$

Matrix of predecessors

MAP Algorithm as Matrix-Vector Product



$$\begin{array}{c}
 v_0 \\
 v_1 \\
 v_2 \\
 v_3
 \end{array}
 \begin{array}{c}
 v_0 \ v_1 \ v_2 \ v_3 \\
 \left(\begin{array}{cccc}
 0 & 0 & 0 & 0 \\
 0 & 0 & 1 & 0 \\
 1 & 1 & 0 & 1 \\
 0 & 0 & 0 & 0
 \end{array} \right)
 \end{array}$$

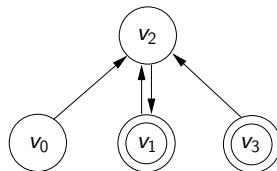
Matrix of predecessors

 \vec{m}
 \perp
 \perp
 \perp
 \perp

Additional data per vertex

\vec{m} – map values

MAP Algorithm as Matrix-Vector Product



$$\begin{array}{c}
 v_0 \ v_1 \ v_2 \ v_3 \\
 v_0 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}
 \end{array}$$

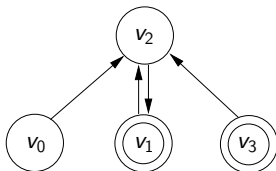
Matrix of predecessors

$$\begin{array}{c}
 \vec{m} \ \vec{o} \\
 \perp \ \perp \\
 \perp \ \perp \\
 \perp \ \perp \\
 \perp \ \perp
 \end{array}$$

Additional data per vertex

\vec{o} – oldmap values

MAP Algorithm as Matrix-Vector Product



$$\begin{array}{c}
 v_0 \ v_1 \ v_2 \ v_3 \\
 v_0 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\
 v_1 \\
 v_2 \\
 v_3
 \end{array}$$

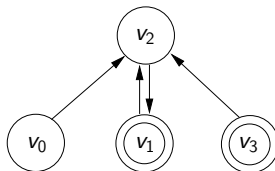
Matrix of predecessors

$$\begin{array}{c}
 \vec{m} \ \vec{o} \ \vec{A} \\
 \perp \ \perp \ 0 \\
 \perp \ \perp \ 1 \\
 \perp \ \perp \ 0 \\
 \perp \ \perp \ 1
 \end{array}$$

Additional data per vertex

\vec{A} – is_accepting

MAP Algorithm as Matrix-Vector Product



$$\begin{array}{c}
 v_0 \ v_1 \ v_2 \ v_3 \\
 v_0 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\
 v_1 \\
 v_2 \\
 v_3
 \end{array}$$

Matrix of predecessors

$$\begin{array}{c}
 \vec{m} \ \vec{o} \ \vec{A} \ \vec{r} \\
 \perp \ \perp \ 0 \ 0 \\
 \perp \ \perp \ 1 \ 0 \\
 \perp \ \perp \ 0 \ 0 \\
 \perp \ \perp \ 1 \ 0
 \end{array}$$

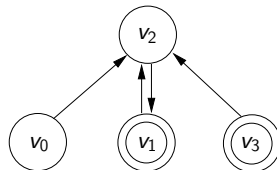
Additional data per vertex

\vec{r} – has been modified

MAP Algorithm as Matrix-Vector Product

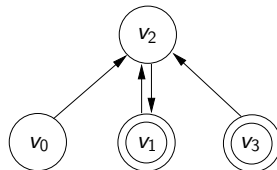
$$\begin{pmatrix} M \end{pmatrix} \times \begin{matrix} V \\ \boxed{} \end{matrix} = \begin{matrix} V' \\ \boxed{} \end{matrix}$$

MAP Algorithm as Matrix-Vector Product – example



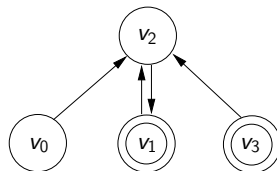
$$\begin{array}{c} v_2 \end{array} \begin{array}{c} v_0 \ v_1 \ v_2 \ v_3 \\ \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{array} \times \begin{array}{ccc} \vec{m} & \vec{o} & \vec{A} \\ \perp & \perp & 0 \\ \perp & \perp & 1 \\ \perp & \perp & 0 \\ \perp & \perp & 1 \end{array} = \begin{array}{cc} \vec{m} & \vec{r} \\ \perp & 0 \\ \perp & 0 \\ \mathbf{3} & \mathbf{1} \\ \perp & 0 \end{array}$$

MAP Algorithm as Matrix-Vector Product – example



$$\begin{matrix} & v_0 & v_1 & v_2 & v_3 & & & & & & \\ v_1 & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \times & \begin{matrix} \vec{m} & \vec{o} & \vec{A} \\ \perp & \perp & 0 \\ \perp & \perp & 1 \\ \mathbf{3} & \perp & 0 \\ \perp & \perp & 1 \end{matrix} & = & \begin{matrix} \vec{m} & \vec{r} \\ \perp & 0 \\ \mathbf{3} & \mathbf{1} \\ 3 & \mathbf{0} \\ \perp & 0 \end{matrix}
 \end{matrix}$$

MAP Algorithm as Matrix-Vector Product – example

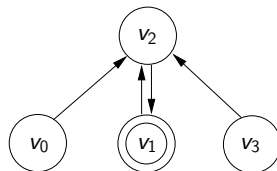


$$\begin{array}{cccc}
 v_0 & v_1 & v_2 & v_3 \\
 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \times & \begin{array}{ccc} \vec{m} & \vec{o} & \vec{A} \\ \perp & \perp & 0 \\ \mathbf{3} & \perp & 1 \\ 3 & \perp & 0 \\ \perp & \perp & 1 \end{array} & = & \begin{array}{cc} \vec{m} & \vec{r} \\ \perp & 0 \\ 3 & \mathbf{0} \\ 3 & 0 \\ \perp & 0 \end{array}
 \end{array}$$

Fix point.

No accepting cycles found

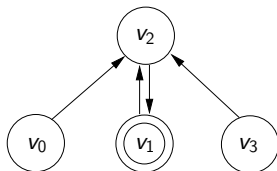
MAP Algorithm as Matrix-Vector Product – example



$$\begin{array}{ccc}
 m & o & \mathcal{A} \\
 \perp & \perp & 0 \\
 3 & \perp & 1 \\
 3 & \perp & 0 \\
 \perp & \perp & 1
 \end{array}
 \implies
 \begin{array}{ccc}
 m & o & \mathcal{A} \\
 \perp & \perp & 0 \\
 3 & \perp & 1 \\
 3 & \perp & 0 \\
 \perp & \perp & \mathbf{0}
 \end{array}$$

Mark v_3 as non-accepting

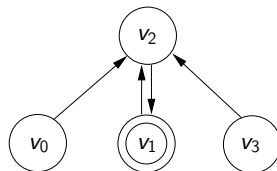
MAP Algorithm as Matrix-Vector Product – example



$$\begin{array}{ccc}
 m & o & \mathcal{A} \\
 \perp & \perp & 0 \\
 3 & \perp & 1 \\
 3 & \perp & 0 \\
 \perp & \perp & 0
 \end{array}
 \implies
 \begin{array}{ccc}
 m & o & \mathcal{A} \\
 \perp & \perp & 0 \\
 \perp & \mathbf{3} & 1 \\
 \perp & \mathbf{3} & 0 \\
 \perp & \perp & 0
 \end{array}$$

Partitioning into subgraphs

MAP Algorithm as Matrix-Vector Product – example



$$\begin{matrix} & v_0 & v_1 & v_2 & v_3 & & & & \\ v_2 & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \times & \begin{matrix} \vec{m} & \vec{o} & \vec{A} \\ \perp & \perp & 0 \\ \perp & 3 & 1 \\ \perp & 3 & 0 \\ \perp & \perp & 0 \end{matrix} & = & \begin{matrix} \vec{m} & \vec{r} \\ \perp & 0 \\ \perp & 0 \\ \mathbf{1} & \mathbf{1} \\ \perp & 0 \end{matrix} \end{matrix}$$

Compact representation

GPU memory is expensive

- Size of the GPU memory is limited
- If a thread accesses too much memory, it blocks other threads

Matrix of predecessors is sparse

- The matrix is made of 1 and 0 only
- Handling full matrices is memory inefficient
- Compact Sparse Row (CSR) representation

Compact representation

GPU memory is expensive

- Size of the GPU memory is limited
- If a thread accesses too much memory, it blocks other threads

Matrix of predecessors is sparse

- The matrix is made of 1 and 0 only
- Handling full matrices is memory inefficient
- Compact Sparse Row (CSR) representation

CSR representation of sparse matrices

$$\begin{array}{c}
 \\
 \\
 1 \\
 2 \\
 3 \\
 4 \\
 5
 \end{array}
 \begin{pmatrix}
 & 1 & 2 & 3 & 4 & 5 \\
 1 & \mathbf{1} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\
 2 & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{1} \\
 3 & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} \\
 4 & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\
 5 & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1}
 \end{pmatrix}$$

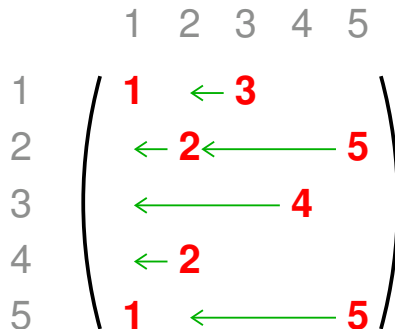
CSR representation of sparse matrices

$$\begin{array}{c}
 \\
 \\
 1 \\
 2 \\
 3 \\
 4 \\
 5
 \end{array}
 \begin{pmatrix}
 & 1 & 2 & 3 & 4 & 5 \\
 \mathbf{1} & & & \mathbf{1} & & \\
 & \mathbf{1} & & & & \mathbf{1} \\
 & & & \mathbf{1} & & \\
 & \mathbf{1} & & & & \\
 \mathbf{1} & & & & & \mathbf{1}
 \end{pmatrix}$$

CSR representation of sparse matrices

$$\begin{array}{c}
 \\
 \\
 1 \\
 2 \\
 3 \\
 4 \\
 5
 \end{array}
 \begin{pmatrix}
 & 1 & 2 & 3 & 4 & 5 \\
 \mathbf{1} & & & \mathbf{3} & & \\
 & \mathbf{2} & & & & \mathbf{5} \\
 & & & & \mathbf{4} & \\
 & \mathbf{2} & & & & \\
 \mathbf{1} & & & & & \mathbf{5}
 \end{pmatrix}$$

CSR representation of sparse matrices



CSR representation of sparse matrices

$$\begin{array}{c}
 \\
 \\
 1 \\
 2 \\
 3 \\
 4 \\
 5
 \end{array}
 \begin{pmatrix}
 & 1 & 2 & 3 & 4 & 5 \\
 \mathbf{1} & \mathbf{3} & & & & \\
 \mathbf{2} & \mathbf{5} & & & & \\
 \mathbf{4} & & & & & \\
 \mathbf{2} & & & & & \\
 \mathbf{1} & \mathbf{5} & & & &
 \end{pmatrix}$$

CSR representation of sparse matrices

1 3
2 5
4
2
1 5

CSR representation of sparse matrices

$M_c = ($ **1** **3** $)$

2 **5**

4

2

1 **5**

CSR representation of sparse matrices

$M_c = (\quad 1 \quad 3 \quad 2 \quad 5 \quad)$



4

2

1 5

CSR representation of sparse matrices

$M_c = (\quad 1 \quad 3 \quad 2 \quad 5 \quad 4 \quad)$



2

1 5

CSR representation of sparse matrices

$M_c = ($
1
3
2
5
4
2
 $)$

↑
↑
↑

1 5

CSR representation of sparse matrices

$M_c = (\quad 1 \quad 3 \quad 2 \quad 5 \quad 4 \quad 2 \quad 1 \quad 5 \quad)$

CSR representation of sparse matrices

$$M_c = (\mathbf{1} \ \mathbf{3} \ \mathbf{2} \ \mathbf{5} \ \mathbf{4} \ \mathbf{2} \ \mathbf{1} \ \mathbf{5})$$

$$M_r = (\mathbf{0} \quad \quad \quad)$$

Diagram illustrating the CSR representation of a sparse matrix. The matrix is shown as two rows: M_c and M_r . The values in M_c are 1, 3, 2, 5, 4, 2, 1, 5, and the value in M_r is 0. Blue arrows point from the values in M_r to the corresponding values in M_c , indicating the mapping of row indices.

CSR representation of sparse matrices

$M_c = (\quad 1 \quad 3 \quad 2 \quad 5 \quad 4 \quad 2 \quad 1 \quad 5 \quad)$



$M_r = (\quad 0 \quad 2 \quad \quad \quad)$

CSR representation of sparse matrices

$M_c = (\quad 1 \quad 3 \quad 2 \quad 5 \quad 4 \quad 2 \quad 1 \quad 5 \quad)$

$M_r = (\quad 0 \quad 2 \quad 4 \quad \quad)$



CSR representation of sparse matrices

$$M_c = (\quad 1 \quad 3 \quad 2 \quad 5 \quad 4 \quad 2 \quad 1 \quad 5 \quad)$$

$$M_r = (\quad 0 \quad 2 \quad 4 \quad 5 \quad 6 \quad)$$

CSR representation of sparse matrices

$$\begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\begin{aligned} \mathbf{Mc} &= (\mathbf{1} \ \mathbf{3} \ \mathbf{2} \ \mathbf{5} \ \mathbf{4} \ \mathbf{2} \ \mathbf{1} \ \mathbf{5}) \\ \mathbf{Mr} &= (\mathbf{0} \ \mathbf{2} \ \mathbf{4} \ \mathbf{5} \ \mathbf{6}) \end{aligned}$$

Memory limitations

GPU Memory consumption

- Data stored on GPU
 - 12B per vertex (4B due to CSR + 8B due to $map, oldmap, \mathcal{A}, r$)
 - 4B per edge (CSR)

$$\begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\begin{aligned} Mc &= (\mathbf{1} \ \mathbf{3} \ \mathbf{2} \ \mathbf{5} \ \mathbf{4} \ \mathbf{2} \ \mathbf{1} \ \mathbf{5}) \\ Mr &= (\mathbf{0} \ \mathbf{2} \ \mathbf{4} \ \mathbf{5} \ \mathbf{6}) \end{aligned}$$

- 1 GB of GPU memory
 - 30 millions of vertices, 150 millions of edges (avg. outdegree 5)
 - 50 millions of vertices, 100 millions of edges (avg. outdegree 2)

Memory limitations

GPU Memory consumption

- Data stored on GPU
 - 12B per vertex (4B due to CSR + 8B due to $map, oldmap, \mathcal{A}, r$)
 - 4B per edge (CSR)

$$\begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\begin{aligned} Mc &= (\color{red}1 \ \color{red}3 \ \color{red}2 \ \color{red}5 \ \color{red}4 \ \color{red}2 \ \color{red}1 \ \color{red}5) \\ Mr &= (\color{blue}0 \ \color{blue}2 \ \color{blue}4 \ \color{blue}5 \ \color{blue}6) \end{aligned}$$

- 1 GB of GPU memory
 - 30 millions of vertices, 150 millions of edges (avg. outdegree 5)
 - 50 millions of vertices, 100 millions of edges (avg. outdegree 2)

Getting the Matrix of Predecessors

Matrix of predecessors

- Must be computed from implicit definition of the graph
- Requires state space generation
- Stored directly in CSR representation
- **CSR representation is constructed on-the-fly**
- Explicit matrix representation is never stored

Matrix of predecessors

Problem

- Successive discovery of predecessors in state space generation
- “Random” writes to the matrix
- Difficult with CSR representation

Graph transposition

- Preserves accepting cycles
- Predecessors are successors in the original graph
- Successors of the original graph are easily generated
- Successive updates to CSR representation

Matrix of predecessors

Problem

- Successive discovery of predecessors in state space generation
- “Random” writes to the matrix
- Difficult with CSR representation

Graph transposition

- Preserves accepting cycles
- Predecessors are successors in the original graph
- Successors of the original graph are easily generated
- Successive updates to CSR representation

Other important features of DiVinE-CUDA

Reduction of the graph stored on GPU

- Exploits automata-based approach to LTL MC
- By analyzing of property automaton some parts of the graph may be known to contain no accepting cycles
- Non-accepting parts are omitted in CSR representation

On-the-fly verification

- CPU and GPU can execute concurrently
- GPU performs MAP algorithm on partially constructed graph
- Accepting cycle is detected prior full state space exploration

Other important features of DiVinE-CUDA

Reduction of the graph stored on GPU

- Exploits automata-based approach to LTL MC
- By analyzing of property automaton some parts of the graph may be known to contain no accepting cycles
- Non-accepting parts are omitted in CSR representation

On-the-fly verification

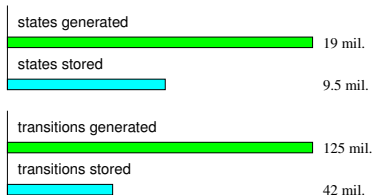
- CPU and GPU can execute concurrently
- GPU performs MAP algorithm on partially constructed graph
- Accepting cycle is detected prior full state space exploration

Impact of graph reduction – Valid properties.

elevator 1



peterson 1



leader

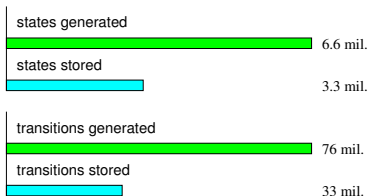


anderson

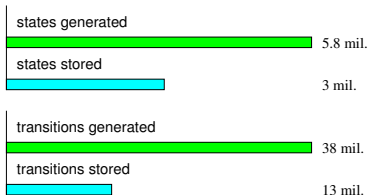


Impact of graph reduction – Invalid properties.

elevator 2



peterson 2



phils



bakery



Comparison to parallel algorithms running on single core

Model	acc. cycle	CUDA MAP			CPU MAP				CPU OWCTY	
		CSR time	CUDA time	total time	1st. iter. time	other iter. time	total time	# iter.	reach. time	total time

elevator 1	N	27	7	34	44	56	100	16	24	41
leader	N	88	1	90	97	600	697	17	90	297
peterson 1	N	107	6	113	175	270	445	16	110	188
anderson	N	32	7	39	64	51	115	5	33	113

elevator 2	Y	34	1	35	50	–	50	1	41	177
phils	Y	47	1	47	295	102	397	5	180	576
peterson 2	Y	26	5	31	173	–	173	1	114	404
bakery	Y	25	1	26	240	–	240	1	219	907

Total time:	386 + 29 = 415	Total time:	2 173	Total time:	2 730
		Speedup:	5.24	Speedup:	6.51

Many-cores vs. Multi-cores (predicted, optimal speedup)

CUDA MAP			CPU MAP	CPU OWCTY
<i>CSR time</i>	<i>CUDA time</i>	<i>total time</i>		

1 core

Total time: $386 + 29 = 415$	Total time: 2 173	Total time: 2 730
Speedup: 5.24	Speedup: 6.51	

2 cores

Total time: $193 + 29 = 222$	Total time: 1087	Total time: 1365
Speedup: 4.87	Speedup: 6.15	

4 cores

Total time: $97 + 29 = 126$	Total time: 544	Total time: 683
Speedup: 4.32	Speedup: 5.42	

8 cores

Total time: $49 + 29 = 78$	Total time: 272	Total time: 342
Speedup: 3.48	Speedup: 4.38	

Conclusion

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration.

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking.

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Pros and Cons

- + On liveness LTL properties outperforms multi-core solutions
- On safety properties does not help
- Can handle verification instances of limited size

Conclusion

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration.

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking.

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Pros and Cons

- + On liveness LTL properties outperforms multi-core solutions
- On safety properties does not help
- Can handle verification instances of limited size

Conclusion

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration.

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking.

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Pros and Cons

- + On liveness LTL properties outperforms multi-core solutions
- On safety properties does not help
- Can handle verification instances of limited size

Conclusion

Acceleration of model checking process by full utilization of modern massively parallel architectures

- successful redesign of Maximal Accepting Predecessor algorithm allowing for significant GPU acceleration.

DiVinE-CUDA

- tool for CUDA Accelerated LTL Model Checking.

Experimental evaluation

- significant reduction of runtimes
- identification of the main bottleneck of the designed approach.

Pros and Cons

- + On liveness LTL properties outperforms multi-core solutions
- On safety properties does not help
- Can handle verification instances of limited size

Future Work

Multi-core preparation of CSR representation

- Ready, will be part of the next release

Ideas to overcome GPU memory limitations

- Partition the matrix of predecessors
- Employ multiple GPU devices

Open problems

- How to CUDA-accelerate state space generation?
- How to CUDA-accelerate CSR representation preparation?

Future Work

Multi-core preparation of CSR representation

- Ready, will be part of the next release

Ideas to overcome GPU memory limitations

- Partition the matrix of predecessors
- Employ multiple GPU devices

Open problems

- How to CUDA-accelerate state space generation?
- How to CUDA-accelerate CSR representation preparation?

Future Work

Multi-core preparation of CSR representation

- Ready, will be part of the next release

Ideas to overcome GPU memory limitations

- Partition the matrix of predecessors
- Employ multiple GPU devices

Open problems

- How to CUDA-accelerate state space generation?
- How to CUDA-accelerate CSR representation preparation?

The End

Thank you four your attention