

Accelerating Domain Propagation: an Efficient GPU-Parallel Algorithm over Sparse Matrices

Boro Sofranac^{1 2} Ambros Gleixner^{1 3} Sebastian Pokutta^{1 2}
{sofranac, gleixner, pokutta}@zib.de

¹Zuse Institute Berlin

²Berlin Institute of Technology

³HTW Berlin

June 25, 2021



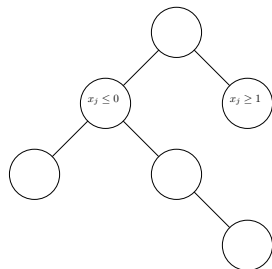
Introduction

Mixed Integer Linear Programming (MIP)

MIP is defined as:

$$\begin{aligned} \min \quad & c^T x \\ \text{s. t.} \quad & Ax \leq b \\ & l_j \leq x_j \leq u_j \quad \forall j \in \mathcal{N} \\ & x_j \in \mathbb{Z} \quad \forall j \in \mathcal{I} \end{aligned}$$

1. Theory: MIPs are \mathcal{NP} -hard
2. Practice: surprisingly fast solvers exist
3. Branch-and-bound algorithm most common



Branch-and-bound search tree

Domain Propagation in MIP

Domain propagation: tightens the bounds of variables

$$-1 \leq x_1 + x_2 + x_3 \leq 1$$

$$-9 \leq x_1 \leq 9$$

$$-1 \leq x_2 \leq 1$$

$$-1 \leq x_3 \leq 1$$

Domain Propagation in MIP

Domain propagation: tightens the bounds of variables

$$-1 \leq x_1 + x_2 + x_3 \leq 1$$

$$-9 \leq x_1 \leq 9$$

$$-1 \leq x_2 \leq 1$$

$$-1 \leq x_3 \leq 1$$

$$u_1^{\text{new}} : x_1 \leq 1 - x_2 - x_3 \leq 3$$

Domain Propagation in MIP

Domain propagation: tightens the bounds of variables

$$-1 \leq x_1 + x_2 + x_3 \leq 1$$

$$-9 \leq x_1 \leq 9$$

$$-1 \leq x_2 \leq 1$$

$$-1 \leq x_3 \leq 1$$

$$u_1^{\text{new}} : x_1 \leq 1 - x_2 - x_3 \leq 3$$

$$l_1^{\text{new}} : x_1 \geq -1 - x_2 - x_3 \geq -3$$

Domain Propagation in MIP

Domain propagation: tightens the bounds of variables

Generalization:

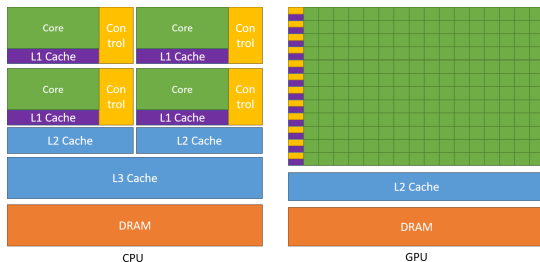
$$\frac{\beta - \bar{\alpha}}{a_j} + u_j \leq x_j \leq \frac{\bar{\beta} - \underline{\alpha}}{a_j} + \ell_j, \quad a_j > 0$$

$$\frac{\bar{\beta} - \underline{\alpha}}{a_j} + u_j \leq x_j \leq \frac{\beta - \bar{\alpha}}{a_j} + \ell_j, \quad a_j < 0$$

Iterated domain propagation

1. Can be seen as a fixed-point iteration
2. May not converge in finite time
3. Tolerance-based termination in practice

MIPs and GPUs



Source: <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

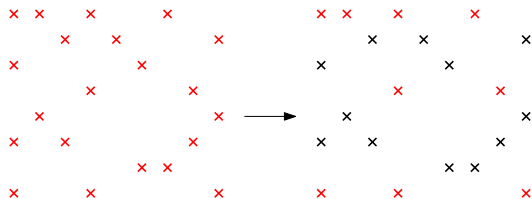
State-of-the-art solvers do not use GPUs. Major challenges:

1. Very irregular and sparse data structures
2. Non-uniform algorithmic behavior

Algorithmic Design

Sequential Propagation Features & GPU Challenges

The constraint marking mechanism



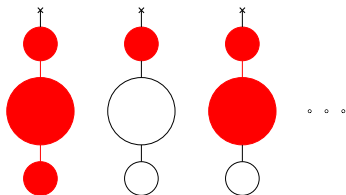
Constraint marking example for two propagation rounds.

GPU implementation challenges:

1. Non-coalesced & random memory accesses
2. Dynamic redistribution of work needed for good load-balancing

Sequential Propagation Features & GPU Challenges

The early termination checks

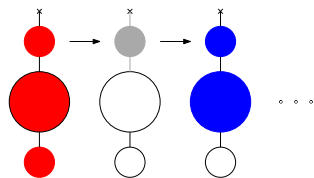


Effect of early termination checks on the workflow of the algorithm

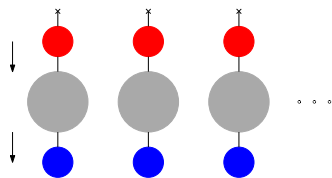
Main GPU implementation challenge:

1. Load-balancing

GPU-parallel Algorithm



Workflow of the sequential algorithm



Workflow of the GPU algorithm

1. Remove features like constraint marking and early termination
2. Switch to throughput-based workflow

Tradeoff: well-structured, static execution at the expense of more work

Handling the Irregular Structure of the Constraint Matrix

Computing min and max activities

Step 1: Compute $\underline{\alpha}$ and $\bar{\alpha}$

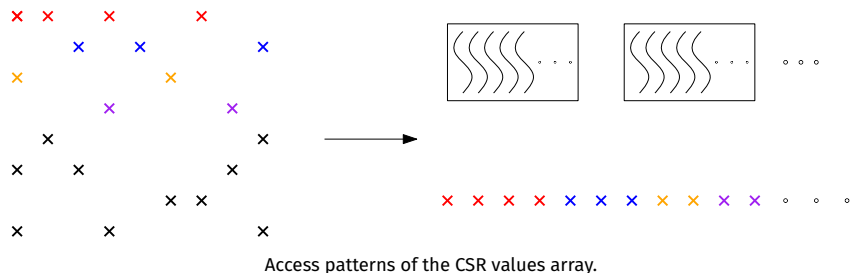
$$\underline{\alpha} := \sum_{i=0}^n a_i b_i \text{ with } b_i = \begin{cases} \ell_i & \text{if } a_i > 0, \\ u_i & \text{if } a_i \leq 0, \end{cases} \quad (1a)$$

$$\bar{\alpha} := \sum_{i=0}^n a_i b_i \text{ with } b_i = \begin{cases} u_i & \text{if } a_i > 0, \\ \ell_i & \text{if } a_i \leq 0, \end{cases} \quad (1b)$$

- Looks similar to SpMV: Ax
- SpMV is highly bandwidth bound
- Memory: Computing $\underline{\alpha}$ and $\bar{\alpha}$ same as computing Al and Au

Approach: Carry over idea from *CSR-adaptive* [Greathouse and Data 2014] and adapt to $\underline{\alpha}$ and $\bar{\alpha}$.

CSR-scalar, CSR-vector, CSR-stream



- *CSR-scalar* exhibits poor memory coalescing
- *CSR-vector* exhibits poor load balancing for short rows
- *CSR-stream*: group rows with few non-zeros together

Computing new bound candidates

- Activities computation: granularity of one thread per non-zero
- $a_j, \ell_j, u_j, \underline{\alpha}_j, \bar{\alpha}_j$ already in shared memory
- atomic updates to handle dependencies

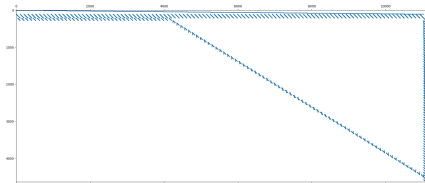
$$\frac{\beta - \bar{\alpha}}{a_j} + u_j \leq x_j \leq \frac{\bar{\beta} - \underline{\alpha}}{a_j} + \ell_j, a_j > 0$$

$$\frac{\bar{\beta} - \underline{\alpha}}{a_j} + u_j \leq x_j \leq \frac{\beta - \bar{\alpha}}{a_j} + \ell_j, a_j < 0$$

Reducing the number of atomic operations

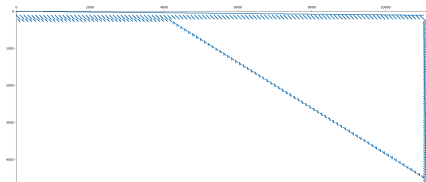
- Bounds are always improving: $\ell^{\text{new}} > \ell^{\text{old}}$
- Therefore, $\ell^{\text{cand}} > \ell^{\text{old}}$, or the bound cannot be updated
- If the bound doesn't improve on the old value, it's not gonna improve on the improvement of the old value!

The Domain Propagation Kernel



Sparsity pattern of MIPLIB instance *drayage-25-27*.

The Domain Propagation Kernel



Sparsity pattern of MIPLIB instance *drayage-25-27*.

- *CSR-stream*, *CSR-vector* adapted to compute $\underline{\alpha}$ and $\bar{\alpha}$
- *CSR-vector*: use one warp
- *CSR-vector-long*: use whole block

Domain Propagation Kernel (1 round)

- 1: **if** end_row - start_row > 1 **then**
 - 2: $\underline{\alpha}$ and $\bar{\alpha}$ via *CSR-stream*
 - 3: **else**
 - 4: **if** nnz < threshold **then**
 - 5: $\underline{\alpha}$ and $\bar{\alpha}$ via *CSR-vector*
 - 6: **else**
 - 7: $\underline{\alpha}$ and $\bar{\alpha}$ via *CSR-vector-long*
 - 8: compute $\ell_{i,j}^{\text{cand}}, u_{i,j}^{\text{cand}}$
 - 9: **if** $u_{i,j}^{\text{cand}} < u_j$ **then**
 - 10: atomic $u_j \leftarrow \min(u_j, u_{i,j}^{\text{cand}})$
 - 11: **if** $\ell_{i,j}^{\text{cand}} > \ell_j$ **then**
 - 12: atomic $\ell_j \leftarrow \max(\ell_j, \ell_{i,j}^{\text{cand}})$
-

Computational Results

Computational Experiments Setup

The test set: MIPLIB 2017

- 1065 MIP instances
- Remove small instances with *variables* < 1000 and *constraints* < 1000
- Divide in 8 subsets of increasing size

Machines:

1. **V100** NVIDIA Tesla V100 PCIe 32GB GPU,
2. **TITAN** NVIDIA Titan RTX 24GB GPU,
3. **RTXsuper** NVIDIA GEFORCE RTX 2080 SUPER 8GB GPU,
4. **P400** NVIDIA Quadro P400 2GB GPU,
5. **amdtr** 64-core AMD Ryzen Threadripper 3990X @ 3.30 GHz with 128 GB RAM CPU,
6. **xeon** 24-core Intel Xeon Gold 6246 @ 3.30GHz with 384 GB RAM CPU,
7. **i7-9700K** 8-core Intel i7-9700K @ 3.60GHz with 64 GB RAM

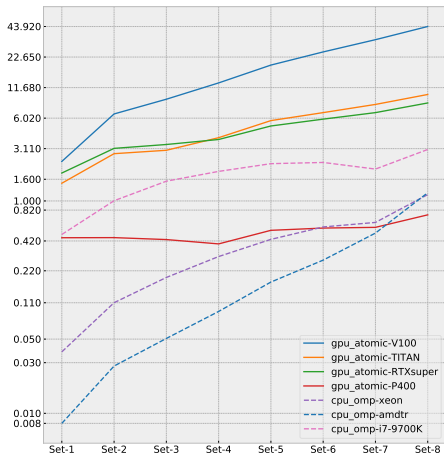
Algorithms:

1. **cpu_seq**
2. **cpu_omp**
3. **gpu_atomic**

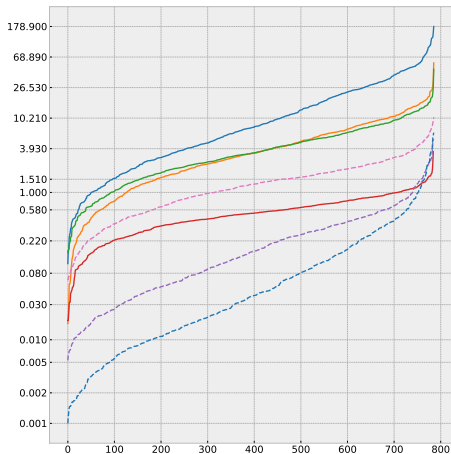
Metric: Speedup over wall-clock time. Base case: **xeon+cpu_seq**

Computational Results: Double precision

(a)



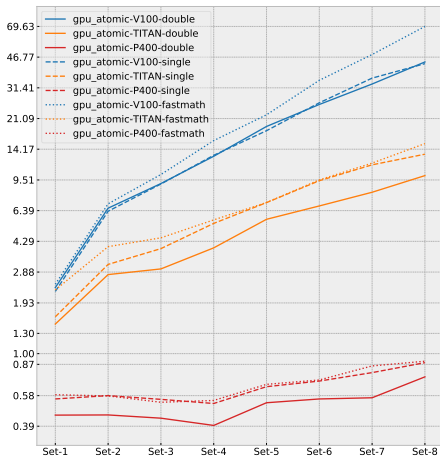
(b)



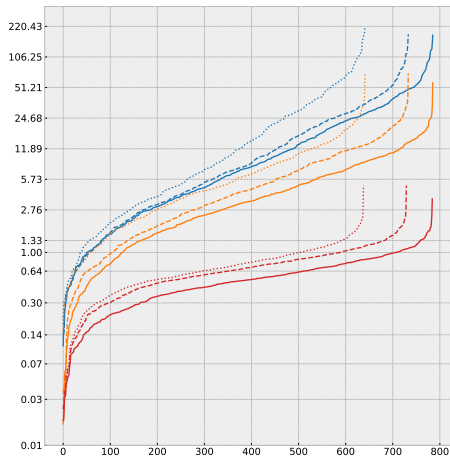
(a) geometric mean of speedups and (b) speedup distributions in ascending order

Computational Results: Single precision

(a)



(b)



(a) geometric mean of speedups and (b) speedup distributions in ascending order