Turbo: Design of a Constraint Programming Solver on GPU

Pierre Talbot pierre.talbot@uni.lu 8th April 2025

University of Luxembourg

UNIVERSITÉ DU LUXEMBOURG

- Declarative paradigm: specify your problem and let the computer solves it for you.
- Many applications: scheduling, bin-packing, hardware design, satellite imaging, ...
- Constraint programming is one approach to solve such combinatorial problems.
- Other approaches include SAT, linear programming, SMT, MILP, ASP,...

Satellite image mosaic



¹Combarro et al., Constraint Model for the Satellite Image Mosaic Selection Problem, CP 2023

Why CP on GPU?

CPU clock speed is stagnating, GPU #cores is increasing quickly each year.



#Cores on Nvidia Tesla cards

Easy speed-up: same code but faster.

"The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin." ^a

^aThe Bitter Lesson, Rich Sutton, 2019,

http://www.incompleteideas.net/IncIdeas/BitterLesson.html

- Machine learning (deep learning, reinforcement learning, ...) has seen tremendous speed-ups (e.g. 100x, 1000x) by using GPU.
- Some (sequential) optimizations on CPU are made irrelevant if we can explore huge state space faster.

Can we replicate the success of GPU on machine learning applications to combinatorial optimization?

Record-Breaking NVIDIA cuOpt Algorithms Deliver Route Optimization Solutions 100x Faster

Mar 20, 2024

🖒 +20 Like 🛛 🤎 Discuss (0)

By Akif Çördük, Piotr Sielski and Moon Chung

Cool, but specialized combinatorial algorithm for routing.

cuOpt: Nvidia Linear Programming Solver



Mittleman's Benchmark cuOpt acceleration compared to state-of-the-art CPU LP (higher is better)

Better! But only for linear constraints over continuous domain.

cuOpt: Nvidia Linear Programming Solver



Better! But only for linear constraints over continuous domain.

Very scarce literature, usually:

- Heuristics: often population-based algorithms².
- Limited set of problems³
- Limited GPU parallelization: offloading to GPU specialized filtering procedures^{4,5}.
- Limited expressivity: solver with max 256 set variables⁶.

For CP-based approach: not general and no proof of correctness.

²A. Arbelaez and P. Codognet, *A GPU Implementation of Parallel Constraint-Based Local Search*, PDP, 2014. ³Jan Gmys. Exactly Solving Hard Permutation Flowshop Scheduling Problems on Peta-Scale GPU-Accelerated Supercomputers. INFORMS Journal on Computing, 2022.

⁴F. Campeotto et al., *Exploring the use of GPUs in constraint solving*, PADL, 2014

⁵F. Tardivo et al., *Constraint propagation on GPU: A case study for the AllDifferent constraint*, Journal of Logic and Computation, 2023.

⁶A. Dovier et al., CUDA: Set Constraints on GPUs, 2022.

A constraint solver on GPU is difficult because:

- Embarassingly parallel (one subproblem per thread) is not efficient on GPU due to limited cache per thread.
- Constraints share variables ⇒ shared-memory parallelism ⇒ synchronization and efficiency issues.

Our Contributions

Theoretical Foundation

We propose a new parallel model of computation:

- Correct: formal proofs of correctness w.r.t. the parallel load/store operations on memory⁷.
- Simple and lock-free: no mutex, no complicated atomic primitives, just barriers.
- Expressive: a process algebra based on *concurrent constraint programming* which supports $\mathbb{Z}, \mathbb{R}, \mathcal{P}(\mathbb{Z}), \ldots$ domains.

⁷P. Talbot et al., A Variant of Concurrent Constraint Programming on GPU, AAAI, 2022.

Our Contributions

Theoretical Foundation

We propose a new parallel model of computation:

- Correct: formal proofs of correctness w.r.t. the parallel load/store operations on memory⁷.
- Simple and lock-free: no mutex, no complicated atomic primitives, just barriers.
- Expressive: a process algebra based on *concurrent constraint programming* which supports $\mathbb{Z}, \mathbb{R}, \mathcal{P}(\mathbb{Z}), \ldots$ domains.

Turbo: First general constraint solver fully executing on GPU.

- General: Support MiniZinc and XCSP3 constraint models (currently only \mathbb{Z} variables).
- Simple: Solving algorithm from 50 years ago.
- Efficient?: Almost on-par with Choco (23% better, 28% worst, 49% equal).
- Open-source: Publicly available on https://github.com/ptal/turbo.

⁷P. Talbot et al., A Variant of Concurrent Constraint Programming on GPU, AAAI, 2022.

On the menu:

- 1. Constraint Programming (CP)
- 2. Towards a **correct** GPU constraint solver.
- 3. Towards an **efficient** GPU constraint solver.

Constraint Programming

Constraint Network

Let X be a finite set of variables and C be a finite set of constraints. A *constraint network* is a pair $P = \langle d, C \rangle$ such that $d \in X \rightarrow Itv$ is the *domain* of the variables where Itv is the set of intervals.

Example

$$\langle \{x \mapsto [0,2], y \mapsto [2,3]\}, \{x \leq y-1\} \rangle$$

A solution is $\{x \mapsto 0, y \mapsto 2\}$.

Interval Propagator

An interval propagator is a function $p_c : \mathbf{D} \to \mathbf{D}$ where $c \in C$ is a constraint and \mathbf{D} the set of all functions $X \to Itv$. Let $d \in \mathbf{D}$, then a propagator is reductive $(p_c(d) \le d)$, monotone $(d \le d' \Rightarrow p_c(d) \le p_c(d'))$ and sound (does not remove solutions).

Example

The propagator for an equality constraint is:

$$p_{x=y}(d) \triangleq d[x \mapsto d(x) \cap d(y), y \mapsto d(x) \cap d(y)]$$

Suppose $d = \{x \mapsto [1,2], y \mapsto [0,5]\}$, then $p_{x=y}(d) = \{x \mapsto [1,2], y \mapsto [1,2]\}$.

```
function SOLVE(d, \{c_1, \dots, c_n\})

d \leftarrow \mathbf{gfp}_d \ p_{c_1} \circ \dots \circ p_{c_n}

if \forall x \in X, \ d(x) = [v, v] then return \{d\}

else if \exists x \in X, \ d(x) = \emptyset then return \{\}

else

\langle d_1, \dots, d_n \rangle \leftarrow \text{split}(d)

return \bigcup_{i=0}^n \text{solve}(d_i, C)
```

end if

end function

- Idea: Divide the problem into many subproblems beforehand (e.g. $N \times 30$ with N the number of threads).
- Intuition: Statistically, there is little chance a subproblem takes longer than the sum of the other subproblems.



⁸A. Malapert et al., 'Embarrassingly Parallel Search in Constraint Programming', JAIR, 2016

Embarrasingly Parallel Search (EPS)

- Idea: Divide the problem into many subproblems beforehand (e.g. $N \times 30$ with N the number of threads).
- Intuition: Statistically, there is little chance a subproblem takes longer than the sum of the other subproblems.



 \Rightarrow **Other approach:** In modern solvers (e.g., Choco, OR-Tools), they use a portfolio approach (e.g., different *split* strategy on the *same problem*).

Embarrasingly Parallel Search (EPS)

- Idea: Divide the problem into many subproblems beforehand (e.g. $N \times 30$ with N the number of threads).
- Intuition: Statistically, there is little chance a subproblem takes longer than the sum of the other subproblems.



On GPU architectures, 1 subproblem per thread is not efficient (limited cache). \Rightarrow Need to parallelize propagation: **gfp** $p_{c_1} \circ \ldots \circ p_{c_n}$.

Towards a Correct GPU Constraint Solver

Concurrent constraint programming (CCP) is a *process calculus* introduced in the eighties⁸. Processes interact by *asking* and *telling* constraints into a *shared grow-only constraint store*.

 $\begin{array}{ll} \langle Program \rangle ::= (t(x_1, \dots, x_n) :- P) + & (process definition) \\ \langle P, Q, \dots \rangle ::= \mbox{when } c \mbox{ then } P & ask \ operator \\ | \ \mbox{tell } c & tell \ operator \\ | \ \ \ P \parallel Q & parallel \ statement \\ | \ \ \exists x, P & hiding \ operator \\ | \ \ \ t(x_1, \dots, x_n) & call \end{array}$

⁸V. A. Saraswat and M. Rinard, *Concurrent constraint programming* (POPL-89)

Concurrent Constraint Programming (CCP)

Concurrent constraint programming (CCP) is a *process calculus* introduced in the eighties⁸. Processes interact by *asking* and *telling* constraints into a *shared grow-only constraint store*.

Example

 $\exists x, y, \text{ (when } x \neq 0 \text{ then tell}(y < x)) \parallel \text{tell}(x > 100)$

A process can work with *partial information*: we do not need to know the exact value of x.

⁸V. A. Saraswat and M. Rinard, *Concurrent constraint programming* (POPL-89)

Concurrent Constraint Programming (CCP)

Concurrent constraint programming (CCP) is a *process calculus* introduced in the eighties⁸. Processes interact by *asking* and *telling* constraints into a *shared grow-only constraint store*.

Clean theoretical framework: every program is a closure operator, unique fixpoint. Can we rely on CCP to implement **gfp** $p_{c_1} \parallel \ldots \parallel p_{c_n}$?

⁸V. A. Saraswat and M. Rinard, *Concurrent constraint programming* (POPL-89)

- Concurrent but not parallel: lack connection with parallel hardware.
- Rely on the (abstract) notion of *constraint system*: good for mathematics, but not for implementation.

 \Rightarrow e.g. **ask** can be an intractable operation depending on the constraint system.

• Unbounded number of variables: parallelism + garbage collection is difficult to achieve efficiently.

We propose a variant of CCP suited for parallel execution on GPU.⁹

- **Parallel**: a process = a thread.
- Correct: equivalence between denotational and (parallel) operational semantics.
- Lattice type: each variable has an explicit domain, instead of being a logical entity.
- No recursion: finite number of variables (known at compile-time).
- Lock-free: processes are executed without locks, even if they share variables.

⁹P. Talbot et al., A Variant of Concurrent Constraint Programming on GPU (AAAI 2022)

Syntax of PCCP

Let $x, y, y_1, \ldots, y_n \in Vars$ be variables, L a lattice, f a monotone function, and b a Boolean variable of type $\langle \{true, false\}, \Leftarrow \rangle$:

$$\begin{array}{l} \langle P, \ Q \rangle ::= \text{ if } b \text{ then } P \\ | & x \leftarrow f(y_1, ..., y_n) \\ | & \exists x: L, \ P \\ | & P \parallel Q \end{array}$$

ask statement tell statement local statement parallel composition

From Constraints to PCCP Programs

Let's x, y be variables, $k \in \mathbb{Z}$ a constant and Itv the interval lattice. We define a function $[\![\varphi]\!]$ compiling a logical formula φ into a PCCP program.

• Constraint $[x + k \le y] \triangleq x \leftarrow f_x(k, x, y) \parallel y \leftarrow f_y(k, x, y)$ with $f_x(k, [x_\ell, x_u], [y_\ell, y_u]) \triangleq [-\infty, y_u - k]$ $f_y(k, [x_\ell, x_u], [y_\ell, y_u]) \triangleq [x_\ell + k, \infty]$

From Constraints to PCCP Programs

Let's x, y be variables, $k \in \mathbb{Z}$ a constant and Itv the interval lattice. We define a function $[\![\varphi]\!]$ compiling a logical formula φ into a PCCP program.

• Constraint $[x + k \le y] \triangleq x \leftarrow f_x(k, x, y) \parallel y \leftarrow f_y(k, x, y)$ with $f_x(k, [x_\ell, x_u], [y_\ell, y_u]) \triangleq [-\infty, y_u - k]$

$$f_{y}(k, [x_{\ell}, x_{u}], [y_{\ell}, y_{u}]) \triangleq [x_{\ell} + k, \infty]$$

- Existential $[\exists x, \varphi] \triangleq \exists x : Itv, [\varphi]$
- Conjunction $\llbracket c_1 \land c_2 \rrbracket \triangleq \llbracket c_1 \rrbracket \parallel \llbracket c_2 \rrbracket$
- Equivalence:

$$\begin{split} \llbracket \varphi \Leftrightarrow \psi \rrbracket &\triangleq \\ & \text{if } entailed(\varphi) \text{ then } \llbracket \psi \rrbracket \\ & \parallel \text{ if } entailed(\psi) \text{ then } \llbracket \varphi \rrbracket \\ & \parallel \text{ if } entailed(\neg \varphi) \text{ then } \llbracket \neg \psi \rrbracket \\ & \parallel \text{ if } entailed(\neg \psi) \text{ then } \llbracket \neg \varphi \rrbracket \end{split}$$

with entailed $(x + k \le y) \triangleq \lceil x \rceil + k \le \lfloor y \rfloor$

Let's consider $[\exists x : Itv, \exists y : Itv, x + 3 \le y \land x + 6 \le z]$

Memory:



Propagators:

$$\begin{vmatrix} \mathbf{x} \leftarrow [-\infty, y_u - 3] \\ y \leftarrow [x_\ell + 3, \infty] \\ \begin{vmatrix} \mathbf{x} \leftarrow [-\infty, z_u - 6] \\ z \leftarrow [x_\ell + 6, \infty] \end{vmatrix}$$

Let's consider $[\exists x : Itv, \exists y : Itv, x + 3 \le y \land x + 6 \le z]$

Memory:

Propagators:



Issue 1: data race? Parallel update of the same integer x_u .

Let's consider $[\exists x : Itv, \exists y : Itv, x + 3 \le y \land x + 6 \le z]$

Memory:

Propagators:



Issue 1: data race? Parallel update of the same integer x_u .

 \Rightarrow **Solution**: Use atomic load and store!

Let's consider $[\exists x : Itv, \exists y : Itv, x + 3 \le y \land x + 6 \le z]$

Memory:

Propagators:



Issue 1: data race? Parallel update of the same integer x_u .

- \Rightarrow **Solution**: Use atomic load and store!
- \Rightarrow In CUDA: Integer load and store are atomic by default!

Let's consider $[\exists x : Itv, \exists y : Itv, x + 3 \le y \land x + 6 \le z]$

Memory:Propagators:x = [1...] $x \leftarrow [-\infty, y_u - 3]$ y = [1..10] $y \leftarrow [x_\ell + 3, \infty]$ z = [1..10] $x \leftarrow [-\infty, z_u - 6]$ $|| \quad z \leftarrow [x_\ell + 6, \infty]$

Issue 1: data race? Parallel update of the same integer x_u .

- \Rightarrow **Solution**: Use atomic load and store!
- \Rightarrow In CUDA: Integer load and store are atomic by default!

Issue 2: nondeterminism? x_u can be equal to 4 or 7 depending on the order of execution.

While something is changing we reexecute all the propagators!



Memory:

Propagators:

fp(

 $\begin{array}{ll} x \leftarrow [-\infty, y_u - 3] \\ || & y \leftarrow [x_\ell + 3, \infty] \\ || & x \leftarrow [-\infty, z_u - 6] \\ || & z \leftarrow [x_\ell + 6, \infty] \end{array}$
While something is changing we reexecute all the propagators!

Memory:



Need to add a condition for progress! We write in the memory only if the value is strictly lower.

 $x \leftarrow v$ iff v < x

- A PCCP process is a reductive and monotone function over a Cartesian product *Store* = L₁ × ... × L_n storing the values of all local variables.
- Since we do not have recursion, we know at compile-time the number of variables.
- Let *Proc* be the set of all PCCP processes.

- A PCCP process is a reductive and monotone function over a Cartesian product *Store* = L₁ × ... × L_n storing the values of all local variables.
- Since we do not have recursion, we know at compile-time the number of variables.
- Let *Proc* be the set of all PCCP processes.

Denotational Semantics

We define a function $\mathcal{D}: Proc \rightarrow (Store \rightarrow Store)$:

$$\begin{aligned} \mathcal{D}(x \leftarrow f(y_1,..,y_n)) &\triangleq \lambda s.s[x \mapsto s(x) \sqcap f(s(y_1),..,s(y_n))] \\ \mathcal{D}(\texttt{if } b \texttt{ then } P) &\triangleq \lambda s.(|s(b) ? \mathcal{D}(P)(s) \circ s |) \\ \mathcal{D}(P \parallel Q) &\triangleq \mathcal{D}(P) \sqcap \mathcal{D}(Q) \end{aligned}$$

Executing the program: **gfp** $\mathcal{D}(P)$.

We obtain the same result if we execute P in parallel or if we replace all parallel \parallel by a sequential operator ; (a transformation we write seq P) defined as follows:

 $\mathcal{D}(P; Q) \triangleq \mathcal{D}(Q) \circ \mathcal{D}(P)$

Let fix f be the set of fixpoints of a function f.

Theorem (Equivalence Between Sequential and Parallel Operators) fix $\mathcal{D}(seq P) = fix \mathcal{D}(P)$ OK Ok, we use a "parallel operator", but it is not really executed in parallel on a machine? Also: $\mathcal{D}(P \parallel Q) \triangleq \mathcal{D}(P) \sqcap \mathcal{D}(Q)$, by unfolding this mathematics definition we have:

$$egin{aligned} \mathcal{D}(P \parallel Q)(S) &= (\mathcal{D}(P) \sqcap \mathcal{D}(Q))(S) \ &= \mathcal{D}(P)(S) \sqcap \mathcal{D}(Q)(S) \end{aligned}$$

What's the problem?

OK Ok, we use a "parallel operator", but it is not really executed in parallel on a machine? Also: $\mathcal{D}(P \parallel Q) \triangleq \mathcal{D}(P) \sqcap \mathcal{D}(Q)$, by unfolding this mathematics definition we have:

$$\mathcal{D}(P \parallel Q)(S) = (\mathcal{D}(P) \sqcap \mathcal{D}(Q))(S) \ = \mathcal{D}(P)(S) \sqcap \mathcal{D}(Q)(S)$$

What's the problem?

We have copied the store!! Not very "Von Neumann Architecture"-friendly. **Contribution:** An operational semantics based on atomic load and store in memory shown equivalent to the denotational semantics under the following assumptions:

(ATOM) Load and store instructions must be atomic.

- (EC) The caches must eventually become coherent.
- (OTA) Values cannot appear *out-of-thin-air*.

Towards an Efficient GPU Constraint Solver

GPU Architecture



5120 cores on a single V100 GPU @ 1290MHz

Whitepaper: https://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf

- Memory coalescence: the way to access the data is important (factor 10).
- **Thread divergence**: each thread within a warp (group of 32 threads) should execute the same instructions.
- **Memory allocation** (dynamic data structures): costly on GPU, everything is generally pre-allocated.
- Other limitations: small cache, limited number of lines of code, limited STL...

Each thread computes its local min (map), then we compute the min of all local min (reduce).

• Map:

3	22	10	23	21	7	91	1	3	10	42	11	8	7	32
Thread 0, $m_0 = 3$				Thread 1, $m_1 = 1$				Thread 2, $m_2 = 3$				Thread 3, $m_3 = 7$		

Iteration 1:

• Map:



Iteration 2:

• Map:



Iteration 3:

• Map:



Iteration 4:

• Map:



3	22	10	23	21	7	91	1	3	10	42	11	8	7	32
T ₀	T_1	T_2	T_3	T_0	T_1	T_2	T_3	T_0	T_1	T_2	T_3	T_0	T_1	T_2

3	22	10	23	21	7	91	1	3	10	42	11	8	7	32
T_0	T_1	T_2	T_3	T_0	T_1	T_2	T_3	T_0	T_1	T_2	T_3	T_0	T_1	T_2





Towards an Efficient GPU Constraint Solver

Representation of Propagators



• Represented using shared_ptr and variant data structures.

 \Rightarrow Uncoalesced memory accesses.

• Code similar to an interpreter:

switch(term.index()) {
 case IVar:
 case INeg:
 case IAdd:
 case IMul:
 // ...

 \Rightarrow Thread divergence.

We simplify the representation of constraints to ternary constraints of the form:

- x = y <op> z where x,y,z are variables.
- The operators are $\{+, /, *, mod, min, max, \leq, =\}$.

Example

The constraint $x + y \neq 2$ is represented by:

t1 = x + yZERO = (t1 = TWO)

where ZERO and TWO are two variables with constant values.

The ternary form of a propagator holds on 16 bytes:

struct bytecode_type {
 int op;
 int x;
 int y;
 int z;
};

- Uniform representation of propagators in memory \Rightarrow coalesced memory accesses.
- Limited number of operators + sorting \Rightarrow reduced thread divergence.

Drawback of TNF: increase in number of propagators and variables.

Benchmark on the MiniZinc Challenge 2024 (89 instances)



The median increase of variables is 4.76x and propagators is 4.33x.

	Bytecode-based + preprocessing	Bytecode-based	Tree-based
Nodes per second	103,036	79,884	14,623
FP iterations per second	1,429,274	1,379,686	88,944
FP iterations per node	13.9	17.3	6.1
Propagators mem. (MB)	0.69	0.91	10.33
Variables mem. (KB)	286.5	397.3	76.6

Analysis of the efficiency of the (preprocessed) bytecode representation of TNF and tree-based representation of arbitrary constraints. The mean is used to aggregate the results. FP stands for fixpoint.

Divergence?



35

Comparison of the best objective values found (timeout: 20 mins, GPU: H100).



What's Next?

Static Scheduling

Find a better ordering of the propagators or fixpoint scheme to reach the greatest fixpoint faster (now, the average is 14 iterations before reaching the fixpoint).

Dynamic Scheduling

- Goal: Avoid executing all propagators in each iteration of the fixpoint loop.
- Constraint network as a sparse graph.
- Waking up propagators using GPU-based algorithms (SpMV and scan).

• Instead of representing the constraints implicitly, we can list all their solutions in a "table":

$$(x \ge 4 \land y > 1 \land z < 3)$$

$$\lor (x = 1 \land y = 2 \land z = 3)$$

$$\lor (x > 1 \land y > 1 \land z > 3)$$

- This representation is useful but has a large and inefficient TNF decomposition.
- \Rightarrow Directly support extension constraints in Turbo.

- Investigate if PCCP (or an extension) can be used for general parallel programming:
 ⇒ finding minimum/maximum in an array, union-find, Floyd-Warshall, ...
- As Turbo is based on lattice theory and abstract interpretation, investigate its applicability to abstract interpretation and verification of neural network.

Conclusion

C++ Abstraction: Lattice Land Project

lattice-land is a collection of libraries abstracting our parallel model.

It provides various data types and fixpoint engine:

- ZLB, ZUB: increasing/decreasing integers.
- B: Boolean lattices.
- VStore: Array (of lattice elements).
- IPC: Arithmetic constraints.
- GaussSeidelIteration: Sequential CPU fixed point loop.
- AsynchronousIteration: GPU-accelerated fixed point loop.

```
• ...
void max(int tid, const int* data, ZLB& m) {
    m.tell(data[tid]);
}
AsynchronousIteration::fixpoint(max);
```

Data races occur rarely, so we should avoid working so much to avoid them.

Properties of the model

A Variant of Concurrent Constraint Programming on GPU (AAAI 2022)¹⁰.

- Correct: Proofs that P; $Q \equiv P || Q$, parallel and sequential versions produce the same results.
- Restartable: Stop the program at any time, and restart on partial data.
- Weak memory consistency: Very few requirements on the underlying memory model ⇒ wide compatibility across hardware, unlock optimization.

¹⁰https://ptal.github.io/papers/aaai2022.pdf

Conclusion: Practical Implementation

"General methods that leverage computation are ultimately the most effective, and by a large margin."—Rich Sutton

Turbo

• Simple: solving algorithms from 50 years ago.

 \Rightarrow no global constraints, nogoods learning, lazy clause generation, restart strategies, event-based propagation, trailing or recomputation-based state restoration and domain consistency.

• Efficient: Almost on-par with Choco (algorithmic optimization VS hardware optimization).

https://github.com/ptal/turbo
Conclusion: Practical Implementation

"General methods that leverage computation are ultimately the most effective, and by a large margin."—Rich Sutton

Turbo

• Simple: solving algorithms from 50 years ago.

 \Rightarrow no global constraints, nogoods learning, lazy clause generation, restart strategies, event-based propagation, trailing or recomputation-based state restoration and domain consistency.

• Efficient: Almost on-par with Choco (algorithmic optimization VS hardware optimization).

https://github.com/ptal/turbo

But...

- Still lagging behind CP+SAT solvers, and SAT learning is inherently sequential...
- There is hope: ACE (a pure CP solver) show CP-only is still competitive in XCSP3 compet.