

A GPU-based Constraint Programming Solver

THE 40TH ANNUAL AAAI CONFERENCE ON ARTIFICIAL INTELLIGENCE
(AAAI 2026)

Pierre Talbot

`pierre.talbot@uni.lu`

<https://ptal.github.io>

24th January 2026

University of Luxembourg



- 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?

State of the Art: Combinatorial Optimization on GPU

Very scarce literature, usually:

- **Heuristics**: often population-based algorithms¹.
- **Limited set of problems**²
- **Limited GPU parallelization**: offloading to GPU specialized filtering procedures^{3,4}.
- **cuOpt**: new MILP solver—relaxation on GPU, search on CPU⁵.

No general-purpose constraint solver on GPU.

¹A. Arbelaez and P. Codognet, *A GPU Implementation of Parallel Constraint-Based Local Search*, PDP, 2014.

²Jan Gmys. Exactly Solving Hard Permutation Flowshop Scheduling Problems on Peta-Scale GPU-Accelerated Supercomputers. *INFORMS Journal on Computing*, 2022.

³F. Campeotto et al., *Exploring the use of GPUs in constraint solving*, PADL, 2014

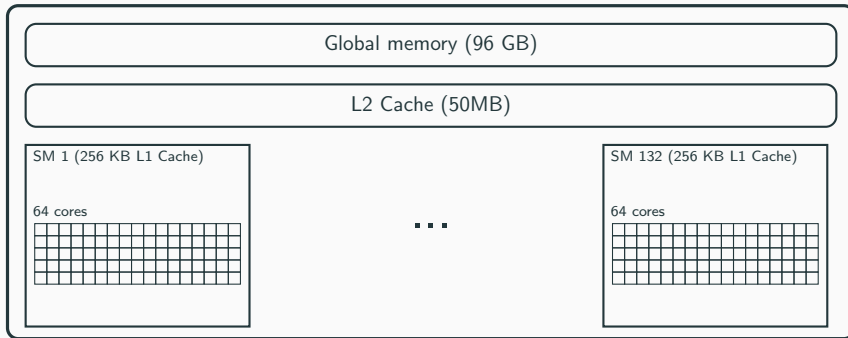
⁴F. Tardivo et al., *Constraint propagation on GPU: A case study for the AllDifferent constraint*, *Journal of Logic and Computation*, 2023.

⁵Using primal-dual linear programming (PDLP).

- **A general constraint solver fully executing on GPU (propagation + search).**
 - ⇒ **General:** Support MiniZinc and XCSP3 constraint models.
 - ⇒ **Simple:** interval-based constraint solving + backtracking search (no global constraints, learning, restart, event-based propagation, ...).
 - ⇒ **Efficient?:** Almost on-par with Choco (23.5% better, 28.6% worst, 48% equal).
 - ⇒ **Open-source:** Publicly available on <https://github.com/ptal/turbo>.
- **Ternary constraint network:** representation of constraints suited for GPU architectures.
- **On-demand subproblems generation strategy:** avoid memory explosion problems by generating problems on-demand.

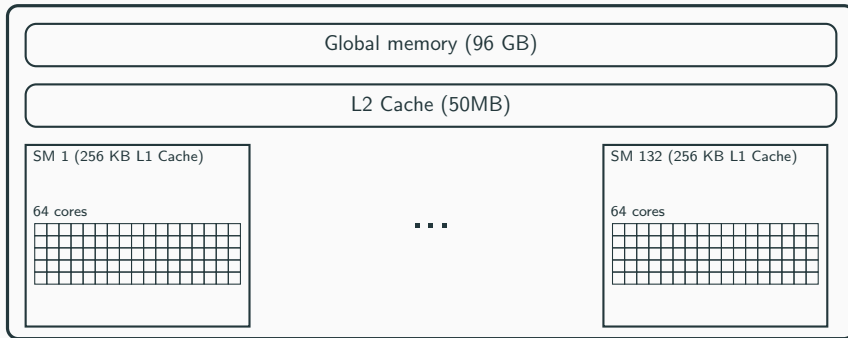
Overview

(Simplified) Architecture of the GPU Nvidia H100



8448 cores grouped in 132 streaming multiprocessors (SM) of 64 cores each.

(Simplified) Architecture of the GPU Nvidia H100



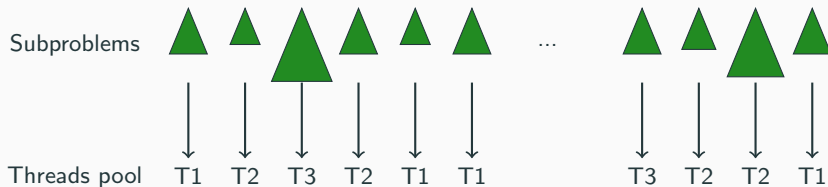
8448 cores grouped in 132 streaming multiprocessors (SM) of 64 cores each.

⇒ **Oversubscribe** (to hide memory latency): 1024 threads per SM

135168 threads running in parallel!

On CPU: Embarrassingly Parallel Search (EPS)⁶

Divide the problem into many subproblems beforehand (e.g. $N \times 30$ with N the number of threads).



⇒ **Other approach:** portfolio approach (e.g., different search strategy on the *same problem*) as seen in Choco and OR-Tools.

Each thread works on its own copy of the problem.

⁶A. Malapert et al., 'Embarrassingly Parallel Search in Constraint Programming', JAIR, 2016

Programming Challenges on GPU

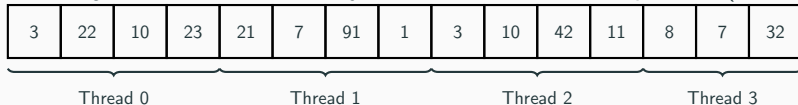
One subproblem per thread is inefficient because:

- **Memory:** suppose each thread needs 1MB (small CSP), then 1GB per SM is constantly moving from global memory to registers.
⇒ Threads competing for cache, slow access to global memory.
- **Single instruction, multiple threads (SIMT):** each consecutive 32 threads should execute the same instructions to avoid thread divergence.
- **Memory coalescence:** the way to access the data is important (factor 10).

Programming Challenges on GPU

One subproblem per thread is inefficient because:

- **Memory:** suppose each thread needs 1MB (small CSP), then 1GB per SM is constantly moving from global memory to registers.
⇒ Threads competing for cache, slow access to global memory.
- **Single instruction, multiple threads (SIMT):** each consecutive 32 threads should execute the same instructions to avoid thread divergence.
- **Memory coalescence:** the way to access the data is important (factor 10).



Programming Challenges on GPU

One subproblem per thread is inefficient because:

- **Memory:** suppose each thread needs 1MB (small CSP), then 1GB per SM is constantly moving from global memory to registers.
⇒ Threads competing for cache, slow access to global memory.
- **Single instruction, multiple threads (SIMT):** each consecutive 32 threads should execute the same instructions to avoid thread divergence.
- **Memory coalescence:** the way to access the data is important (factor 10).

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

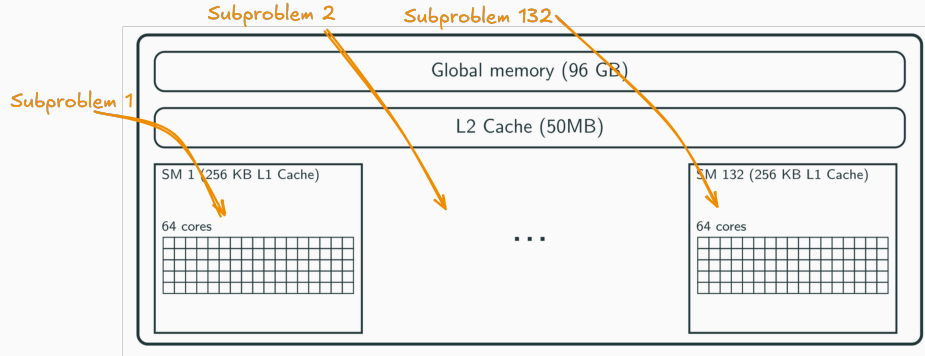
Programming Challenges on GPU

One subproblem per thread is inefficient because:

- **Memory:** suppose each thread needs 1MB (small CSP), then 1GB per SM is constantly moving from global memory to registers.
⇒ Threads competing for cache, slow access to global memory.
- **Single instruction, multiple threads (SIMT):** each consecutive 32 threads should execute the same instructions to avoid thread divergence.
- **Memory coalescence:** the way to access the data is important (factor 10).

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

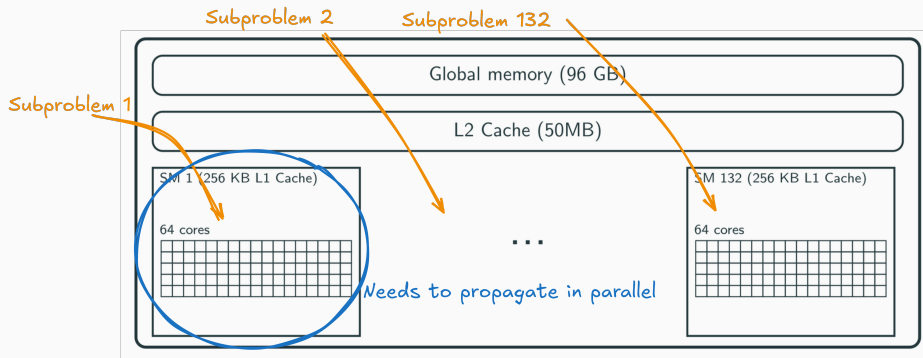
One subproblem per SM⁷ with EPS.



Less memory transfer, L1 cache per subproblem.

⁷More precisely, one subproblem per block, a block is running on a single SM. Several blocks can be scheduled on the same SM.

Parallel Propagation

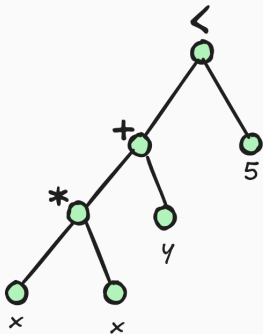


We previously proposed a correct model of lock-free parallel propagation, but lacked efficiency⁸.

⁸P. Talbot et al., *A Variant of Concurrent Constraint Programming on GPU*, AAAI, 2022.

Ternary Constraint Network

Representation of Propagators



- Represented using `shared_ptr` and variant data structures.

⇒ Uncoalesced memory accesses.

- Code similar to an interpreter:

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

⇒ Thread divergence.

Ternary Constraint Network

CSP $\langle X, D, \{c_1, \dots, c_n\} \rangle$ where each c_i is of the form $x = y \text{ <op> } z$ with:

- $x, y, z \in X$ (no constant),
- $op \in \{+, /, *, \textit{mod}, \textit{min}, \textit{max}, \leq, =\}$.

Expressive enough to support all problems of MiniZinc competitions 2022–2024.

Ternary Constraint Network

CSP $\langle X, D, \{c_1, \dots, c_n\} \rangle$ where each c_i is of the form $x = y \text{ <op> } z$ with:

- $x, y, z \in X$ (no constant),
- $op \in \{+, /, *, \textit{mod}, \textit{min}, \textit{max}, \leq, =\}$.

Expressive enough to support all problems of MiniZinc competitions 2022–2024.

Example

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

$x = t1 + y$

$\text{ZERO} = (t1 = \text{TWO})$

equivalent to $t1 = x - y$

equivalent to $\text{false} \Leftrightarrow (t1 = 2)$

where ZERO and TWO are two variables with constant values.

Bytecode Representation

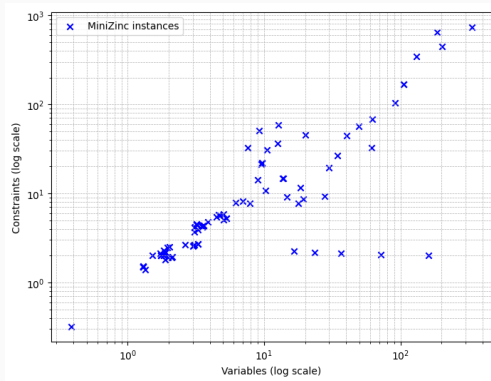
The ternary form of a propagator holds on 16 bytes:

```
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 TCN: increase in number of propagators and variables.

Benchmark on the MiniZinc Challenge 2024 (96 instances)⁹.



The **median increase** of variables is 4.45x and propagators is 4.34x.

The **maximum increase** of variables is 336x and propagators is 731x.

⁹Instances not solved during preprocessing.

Benchmarking

Experimental Evaluation

On 98 instances of the MiniZinc 2024 competition.¹⁰

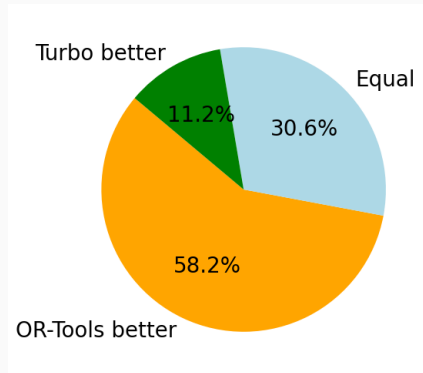
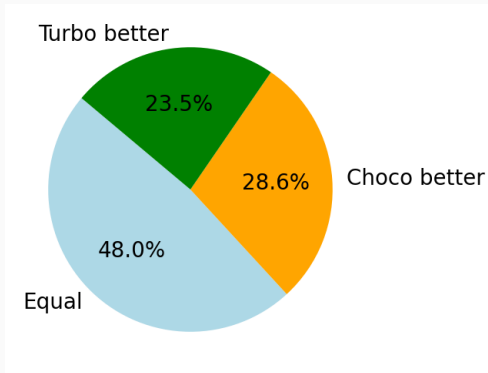
Timeout 20 minutes, CPU 64 cores, GPU H100.

solver	MiniZinc score	#Optimal
Or-Tools 9.9 (64 threads)	266.7	82
Choco 4.10.18 (64 threads)	190.8	44
Or-Tools 9.9 (fixed search)	119.9	37
Choco 4.10.18 (fixed search)	49.3	25
Turbo 1.2.8 (fixed search)	45.8	20

¹⁰1 instance unsat at root, 1 instance for which TCN is too large.

1-to-1 Comparison

Comparison of the best objective values found.



Conclusion

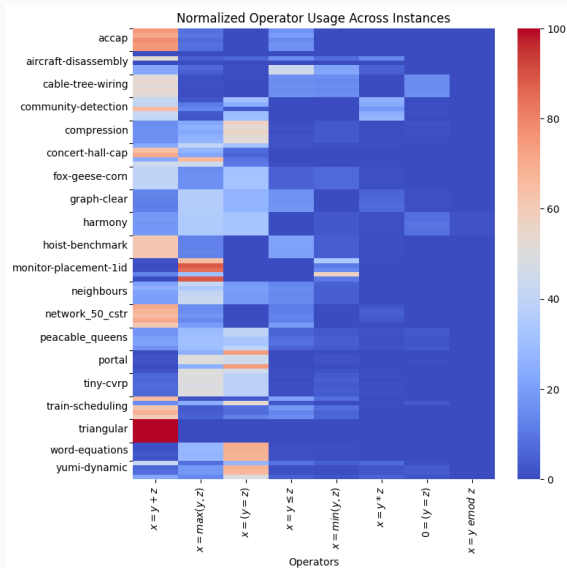
Turbo: General-purpose GPU constraint solver

- **Simple:** solving algorithms from 50 years ago.
⇒ 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).
- Many possible optimizations to improve the efficiency, but need to be redesigned for GPU.



<https://github.com/ptal/turbo>

Divergence?



Lock-free Parallel Propagation

Example of Parallel Propagation¹¹

Let's consider $\mathcal{I}[\![x \leq 4 \wedge x \leq 5]\!] = \mathcal{I}[\![x \leq 4]\!] \parallel \mathcal{I}[\![x \leq 5]\!]$

Memory:

$$x = [-\infty, \infty]$$

Propagators:

$$\begin{array}{ll} x \leftarrow [-\infty, 4] & (\mathcal{I}[\![x \leq 4]\!]) \\ \parallel & \\ x \leftarrow [-\infty, 5] & (\mathcal{I}[\![x \leq 5]\!]) \end{array}$$

¹¹P. Talbot et al., *A Variant of Concurrent Constraint Programming on GPU*, AAAI, 2022.

Example of Parallel Propagation¹¹

Let's consider $\mathcal{I}[\![x \leq 4 \wedge x \leq 5]\!] = \mathcal{I}[\![x \leq 4]\!] \parallel \mathcal{I}[\![x \leq 5]\!]$

Memory:

$x = [-\infty, ?]$

Propagators:

$x \leftarrow [-\infty, 4] \quad (\mathcal{I}[\![x \leq 4]\!])$
 $\parallel \quad x \leftarrow [-\infty, 5] \quad (\mathcal{I}[\![x \leq 5]\!])$

Issue: nondeterminism? x can be equal to $[-\infty, 4]$ or $[-\infty, 5]$ depending on the order of execution.

¹¹P. Talbot et al., *A Variant of Concurrent Constraint Programming on GPU*, AAAI, 2022.

Example of Parallel Propagation¹¹

Let's consider $\mathcal{I}[x \leq 4 \wedge x \leq 5] = \mathcal{I}[x \leq 4] \parallel \mathcal{I}[x \leq 5]$

Memory:

$x = [-\infty, 4]$

Propagators:

$x \leftarrow [-\infty, 4] \quad (\mathcal{I}[x \leq 4])$
 $\parallel \quad x \leftarrow [-\infty, 5] \quad (\mathcal{I}[x \leq 5])$

Issue: nondeterminism? x can be equal to $[-\infty, 4]$ or $[-\infty, 5]$ depending on the order of execution.

\Rightarrow **Solution:** fixpoint + fair scheduling + strict updates (if $(v < x.\text{ub}) \{ x.\text{ub} = v; \}$).

¹¹P. Talbot et al., *A Variant of Concurrent Constraint Programming on GPU*, AAAI, 2022.

GPU Fixpoint Algorithm

```
__device__ void fixpoint(Store& d, Props* props, int n) {  
    __shared__ bool has_changed = true;  
    // Keep going until no variable domain is modified.  
    while(has_changed) {  
        __syncthreads(); has_changed = false; __syncthreads();  
    }
```

GPU Fixpoint Algorithm

```
__device__ void fixpoint(Store& d, Props* props, int n) {  
    __shared__ bool has_changed = true;  
    // Keep going until no variable domain is modified.  
    while(has_changed) {  
        __syncthreads(); has_changed = false; __syncthreads();  
        // Execute all propagators (similar to AC1)  
        for(int i = threadIdx.x; i < n; i += blockDim.x) {  
            has_changed |= props[i].propagate(d);  
        }  
        __syncthreads();  
    }  
}
```

GPU Challenges

- Coalesced memory accesses of the propagator representation `props[i]`.
- Avoiding divergence in `propagate`.