Towards Chapel-based Exascale Tree Search Algorithms: dealing with multiple GPU accelerators

Tiago Carneiro, Nouredine Melab, Akihiro Hayashi, Vivek Sarkar

University of Luxembourg (Luxembourg), INRIA Lille/Université de Lille (France), Georgia Tech – (USA)

> March 22–27, 2021 (virtual event)





Agenda

Introduction

Initial Premises

The Proposed Algorithm

Performance evaluation

Discussion/Conclusion

Algorithms for solving combinatorial optimization problems (COPs) can be divided into exact (complete) or approximate strategies.

Approximate: feasible time.

- Local Search (k-opt)
- Meta-heuristics (GRASP, Tabu Search)
- Hybrid Methods, etc.

Exact: exponential execution time.

- Tree-based search: backtracking, B&B (and their variations).
- Dynamic Programming, etc.

Tree Search Algorithms <u>implicitly</u> enumerate a solution space, dynamically building a tree.



- Nodes are removed from the active set according to the search strategy.
- The search continues until the active set is empty.
- (Crainic, Le Cun, Roucairol (2006)).

Tree search algorithms...

- are highly irregular,
- demand hand-optimized data structures,
- demand load balancing schemes,
- exponential execution time.

Implementation is

performance-oriented :

- $\, \bullet \,$ Usually in C/C++
- Low-level features \rightarrow performance
- Parallel computing libraries

GPUs are crucial in exact optimization : they enable solving to the optimality instances having prohibitive execution times on CPUs.

In this way, it is expected that exascale systems will decrease the time required to solve instances of COPs to the optimality.

 several programming models, languages, runtimes combined to efficiently exploit all levels of parallelism of such heterogeneous systems.

High-productivity languages:

- Exascale systems are going to be complex to program.
- Efforts towards productivity are crucial for better exploiting the future generation of supercomputers.

Chapel plays a special role: it provides different parallel and distributed iterators that implement load balancing among processes, besides other high-level features.

Chapel's parallel iterators:

- Iterators in Chapel are similar to procedures that can be used to isolate iterations from the loop body.
- Each value yielded by the iterator corresponds to an iteration of the loop.

Chapel distributed iterators are a key feature for achieving a trade-off between productivity and parallel efficiency/performance in distributed tree-based optimization.

Distributed iterators implement:

- Master/worker model for pool management;
- Distributed load balancing/work distribution;
- Complex distributed metrics reduction;
- Termination criteria.

Thus, using Chapel can be up to $8 \times$ more productive than MPI+C in the context of this work (Carneiro et al. (2020)).

Agenda

Introduction

Initial Premises

The Proposed Algorithm

Performance evaluation

Discussion/Conclusion

Initial Premises

GPUs are required in exact optimization : they enable solving to the optimality instances having prohibitive execution times on CPUs.

GPU support in Chapel:

- A couple of approaches try to mitigate the gap between Chapel and GPUs (Chu et al. (2017) and Sidelnik et al. (2012)).
- They do not allow lower-level GPU programming features.
- We do need such features in exact optimization.

Initial Premises

Hayashi et al. (2019): https://github.com/ahayashi/chapel-gpu

- GPU/CPU oriented iterator for Chapel.
- ${\circ}$ Allows the use of pre-compiled CUDA-C/C++ kernels.
- Allows concurrent heterogeneous distributed CPU/GPU execution.
- However, does not provide load balancing between GPUs/CPU tasks.

Tree-based search are highly irregular applications Load balancing is therefore crucial.

Initial Premises

Hayashi's iterator is not suitable for irregular GPU-based applications.

So, how to use both Chapel's high-level features for distributed programming and GPUs?

Agenda

Introduction

Initial Premises

The Proposed Algorithm

Performance evaluation

Discussion/Conclusion

The Proposed Algorithm

We revisit:

- Master-worker distributed tree-search,
- Using Chapel's distributed iterators for work distribution,
- Distributed backtracking,
- Enumerates *all* valid configurations of the N-Queens.

But now:

- Using C-Interoperability layer along with Chapel's high-level features,
- Pre-compiled CUDA-C kernels,
- CPU-GPU scheme.

the N-Queens problem is a proof-of-concept that motivates further improvements in solving related COPs.





- On Locale 0 (master) -Task 0.
- Serial search.



- On Locale 0 (master) -Task 0.
- Serial search.
- Starts with the initial configuration of the problem.



 Partial search until a cutoff depth.

Generating the initial pool of subproblems: partial search on CPU.



 Pool keeps all valid and feasible partial solutions with cutoff elements of the permutation.

Master-worker: distributed load balancing using iterators.

Algorithm 1: The Master locale.

 $1 N \leftarrow qet_problem()$ $2 \ cutoff \leftarrow aet \ cutoff \ depth()$ 3 second cutoff \leftarrow get scnd cutoff depth() 4 $P \leftarrow \{\}$ Node 5 metrics \leftarrow (0,0) 6 $metrics + = initial_search(N, cutof f, P)$ 7 Size $\leftarrow \{0..(|P|-1)\}$ // Domain **8** $D \leftarrow Size$ mapped onto locales to a standard distribution 9 $P_d \leftarrow [D]$: Node 10 $P_d = P //$ Using implicit bulk-transfer 11 forall node in P_d following a distributed iterator with (+ reduce metrics) do $metrics + = Algorithm_2(N, node, cutoff,$ 12 second cutoff) 13 14 end

15 present_results(metrics)

Master-worker: distributed load balancing using iterators.



From one node received via parallel iterator, how to...

- Generate load enough for all GPUs of the system?
- How to divide the load among GPUs?

From one node received via parallel iterator, how to...

- 1. Generate load enough for all GPUs of the system?
 - Nested parallelism generating a local pool via partial search.
 - From depth **cutoff** until depth **second_cutoff**.
 - Thus, it is another partial search on CPU.
 - Task-local pool.
- 2. How to divide the subproblems of the local pool among GPUs?

From one node received via parallel iterator, how to...

- 1. Generate load enough for all GPUs of the system?
 - **Nested parallelism** for generating a local pool via partial search.
 - From depth cutoff until depth second_cutoff.
 - Thus, it is another partial search on CPU.
 - Task-local pool.
- 2. From the task-local pool...
 - How to divide the subproblems of the local pool among GPUs?

How to divide the local pool among GPUs?

Algorithm 3: Exploiting multiple GPUs.	
Input: N, P, the second cutoff depthOutput: A tuple containing the explored tree size and the number of complete and valid solutions found on GPU.	
$1 tree_h \leftarrow [0 P - 1] int$ $2 sols_h \leftarrow [0 P - 1] int$ $3 \gamma \leftarrow cuda_get_num_devices()$ $4 GPU load \leftarrow P $ $5 forall gpu_id in 0 \gamma - 1 do$ $6 \qquad cuda_set_amu(amu id)$	 Variables for metrics reduction.
$\begin{array}{llllllllllllllllllllllllllllllllllll$	 Getting the number of GPUs via
12 call_GPU_search(N, depth, device_load, pool_ptr, 13 tree_ptr, sols_ptr) 14 end	CUDA.
$\begin{array}{llllllllllllllllllllllllllllllllllll$	
18 return (redTree, redSols)	

How to divide the local pool among GPUs?

Algorithm 3: Exploiting multiple GPUs.	
Input: N, P, the second cutoff depth Output: A tuple containing the explored tree size and the number of complete and valid solutions found on GPU.	
1 tree_h $\leftarrow [0 P - 1]$ int 2 sols_h $\leftarrow [0 P - 1]$ int 3 $\gamma \leftarrow cuda_get_num_devices()$	
$\begin{array}{c c} \textbf{4} & GPU_load \ \leftarrow \ P \\ \textbf{5} & \textbf{forall} & gpu_id \ in \ 0\gamma - 1 \ \textbf{do} \\ \hline \textbf{6} & \hline \begin{array}{c} cuda_set_gpu(gpu_id) \\ device_load \ \leftarrow \ get_load(gpu_id, GPU_load, \gamma) \\ \textbf{s} & \\ \textbf{sols_ptr} \ \leftarrow \ sols_h \ + \ stride \\ \hline \textbf{0} & \\ \textbf{tree_ptr} \ \leftarrow \ tree_h \ + \ stride \\ \end{array} \right)$	 For each GPU Get the GPU load.
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	
$\begin{array}{ll} reasons \leftarrow (+ reauce sois_n) \\ retrics+ = (redTree, redSol) \\ return (redTree, redSols) \end{array}$	

How to divide the local pool among GPUs?

Algorithm 3: Exploiting multiple GPUs. **Input:** N, P, the second cutoff depth Output: A tuple containing the explored tree size and the number of complete and valid solutions found on GPU. 1 tree_h $\leftarrow [0..|P|-1]$ int 2 sols $h \leftarrow [0..|P| - 1]$ int $3 \gamma \leftarrow cuda \ get \ num \ devices()$ 4 GPU_load \leftarrow |P| 5 forall qpu_id in $0..\gamma - 1$ do cuda set qpu(qpu id) 6 device_load \leftarrow get_load(gpu_id, GPU_load, γ) 7 $stride \leftarrow qet_starting_point(GPU_load, qpu_id, \gamma)$ 8 sols $ptr \leftarrow sols h + stride$ 9 tree $ptr \leftarrow tree h + stride$ 10 $pool_ptr \leftarrow P + stride$ 11 call_GPU_search(N, depth, device_load, pool_ptr, 12 tree ptr. sols ptr) 13 14 end 15 $redTree \leftarrow (+ reduce tree_size_h)$ 16 redSols \leftarrow (+ reduce sols h) 17 metrics + = (redTree, redSol)18 return (redTree. redSols)

- For each
 GPU...
- Calculate its load.
- Calculate the pointers for calling the CUDA kernel.
- Strides on the memory.

How to divide the local pool among GPUs?

Algorithm 3: Exploiting multiple GPUs.	
Input: N, P, the second cutoff <i>depth</i> Output: A <i>tuple</i> containing the explored tree size and the number of complete and valid solutions found on GPU.	
1 tree_h $\leftarrow [0 P - 1]$ int 2 sols_h $\leftarrow [0 P - 1]$ int 3 $\gamma \leftarrow cuda_get_num_devices()$ 4 GPU_load $\leftarrow P $	• Paralle
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	using Chape high-le feature
 14 end 15 redTree ← (+ reduce tree_size_h) 16 redSols ← (+ reduce sols_h) 17 metrics+ = (redTree, redSol) 18 return (redTree, redSols) 	

el ion el's evel es.

Overview of the algorithm

Two levels of parallelism: intra-node and inter-node.



Agenda

Introduction

Initial Premises

The Proposed Algorithm

Performance evaluation

Discussion/Conclusion

Evaluation: The following programs were conceived for enumerating all valid and complete configurations of the N-Queens problem.

- **Baseline**: single-node multi-GPU implementation optimized for single-locale execution written in CUDA-C.
- **GPUIterator**: distributed version of the baseline implementation written in Chapel. Uses Hayashi et al. (2019). No load balancing.
- ChpIGPU: implementation previously detailed.
- Obs. : All implementations employ the same CUDA-C kernel code.

Parameters settings:

- 12 computer nodes.
- *Two* Intel Xeon E5-2650 v4@ 2.00GHz (a total of 24 cores/ 48 threads per node) and 128 GB RAM.
- *Two* NVIDIA GeForce GTX 1080 Ti Pascal generation (11GB RAM and 3584 CUDA cores @ 1582Mhz).
- Maximum 24 GPUs (86,016 CUDA cores) used in the experiments.
- 100 Gbps Intel Omni-Path network.

Single-node performance: one computer node - two GPUs.



Single-node performance: one computer node – two GPUs.



- Chpl-based GPUIterator equivalent to the CUDA-C baseline.
- Using the iterator does not mean significant overhead.

Single-node performance: *one* computer node – *two* GPUs.



- The ChpIGPU implementation is not single-node oriented.
- Nested parallelism and other features for distributed execution.

Single-node performance: one computer node - two GPUs.



The overhead becomes less significant as N increases.
 It is from

 It is from 2.15× to 10% slower than the baseline.

Distributed performance: for 4 GPUs (2 computer nodes) to 24 GPUs (12 computer nodes).



 For the smallest instance, that take only a few seconds, GPUIterator is faster than ChpIGPU.

Distributed performance: for 4 GPUs (2 computer nodes) to 24 GPUs (12 computer nodes).



 For the smallest instance, that take only a few seconds, GPUIterator is faster than ChpIGPU.

 As more GPUs (nodes) are added...

Distributed performance: for 4 GPUs (2 computer nodes) to 24 GPUs (12 computer nodes).



 GPUlterator performs poorly due to the lack of distributed load balancing.

Distributed performance: for 4 GPUs (2 computer nodes) to 24 GPUs (12 computer nodes).



• For N = 18and 16–24 GPUs, ChpIGPU is from $1.32 \times -1.45 \times$ faster.

Distributed performance: for 4 GPUs (2 computer nodes) to 24 GPUs (12 computer nodes).



 For N ≥ 20, ChplGPU is from 1.13× (4 GPUs) to 1.77× (24 GPUs) faster than its counterpart.

Efficiency: speedup compared to the CUDA-C baseline.



 As a consequence... For N ≥ 18, the speedups achieved by the GPUIterator implementation are around 50% of the linear speedup.

Carneiro et al.

HPCS 2020

Efficiency: speedup compared to the CUDA-C baseline.



• In turn, for the ChpIGPU implementation, when $N \ge 19$, the speedups achieved for all < problem, #GPUS > configurations are at least 80% of the linear speedup. For $N \ge 20$ ranged from 87% to 91%.

Carneiro et al.

HPCS 2020

Agenda

Introduction

Initial Premises

The Proposed Algorithm

Performance evaluation

Discussion/Conclusion

Discussion

GPU Iterator:

- The best option for a single-locale implementation.
- Best option in terms of time to a first implementation.
- Low programming effort to get a distributed version.

Error-prone details are hidden:

- Pointer arithmetics
- Load distribution (CPU-GPU, GPUs, and locales)
- Small SLOC count (8 lines+)

However, poor scalability in <u>irregular</u> tree search. Good solution for programming distributed heterogeneous and <u>regular</u> applications (*).

Discussion

ChpIGPU:

- Nested parallelism.
- GPU-related load distribution by hand.
- Much more complex.

Error-prone details explicitly programmed:

- $\bullet~1.5\times$ longer than using the iterator.
- GPU load distribution.
- Pointer arithmetics.

Pays off : scales much better than its counterpart.

Discussion

However, some challenges concerning the use of GPUS remain, even using a high-productivity language:

Challenges:

- CPU-GPU heterogeneity.
- How to generate on CPU load enough for all CPU/GPUs?
- How to perform load balancing between all CPU/GPUs using iterators?

First research direction: incorporate work-stealing mechanism into Hayashi's GPUIterator module.

Future Research Directions

Research directions:

- Extend the current implementation to a distributed B&B.
- Solving challenging COPs (FSP, Q3AP, etc.).

Road Towards Exascale:

- The scalability should be increased
- CPU-GPU heterogeneity?
- Fault tolerance. Checkpointing?

Thank you!

Questions ?

https://github.com/tcarneirop/ChOp https://github.com/ahayashi/chapel-gpu

tiago.carneiropessoa@uni.lu

