

Electronics and Computer Science
Faculty of Physical and Applied Sciences
University of Southampton

Kim Svensson
December 11, 2012

Combinatorial optimization of NP-hard problems using
dynamic GPU programming

Project supervisor: Dr. Sarvapali D. Ramchurn
Second examiner: Dr. Sarvapali D. Ramchurn

A project progress report submitted for the award of
CompSci MEng

Abstract

Solving problems that have an exponential growth in time and space have always been a difficulty as the advancement in computing speed only allows for a small increase of the problem size. The advent of parallel computing using a GPU have enabled a speedup when it comes to problems which have some data and computation independent parts. The purpose of this project is to develop and evaluate dynamic programming algorithms suitable for parallel computation to find an optimal solution in the domain of complete coalition formation and combinatorial auctions. The heterogeneous algorithm is run using the Nvidia CUDA framework and preliminary result shows a substantial difference and is publishable results.

Contents

- 0.1 Background 1
 - 0.1.1 Combinatorial Auctions and Coalition Structure Generation 1
 - 0.1.2 NVIDIA CUDA 2
- 0.2 Literature review 3
- 0.3 Work 4
 - 0.3.1 Result 7
 - 0.3.2 Evaluation 7
- 0.4 Future work 8

0.1 Background

0.1.1 Combinatorial Auctions and Coalition Structure Generation

A combinatorial auction can be described as an autioneer having n assets to sell where bidders may bid on sets of the assets. The goal is to calculate which set of disjoint bids gives the most revenue to the auctioneer. Consider that the auctioneer is selling the assets \clubsuit , \diamond and \heartsuit . The bids collected are as following:

Set	Bid
\clubsuit	3
\diamond	2
\heartsuit	6
\clubsuit, \diamond	7
\clubsuit, \heartsuit	7
\diamond, \heartsuit	4
$\clubsuit, \diamond, \heartsuit$	6

Table 1: This table shows some data

The winning bid configuration would be the following set: $\{\{\heartsuit\}, \{\clubsuit, \diamond\}\}$ with a total sum of $7\{\clubsuit, \diamond\} + 6\{\heartsuit\} = 13$ as it is the maximum sum of all combinations of disjoint sets. Where all complete disjoint sets are:

Complete disjoint set	Sum
$\{\{\heartsuit\}, \{\clubsuit, \diamond\}\}$	13
$\{\{\diamond\}, \{\clubsuit, \heartsuit\}\}$	9
$\{\{\clubsuit\}, \{\diamond, \heartsuit\}\}$	7
$\{\{\clubsuit, \diamond, \heartsuit\}\}$	6
$\{\{\clubsuit\}, \{\diamond\}, \{\heartsuit\}\}$	11

Table 2: This table shows some data

The way you calculate the optimal solution for combinatorial auctions is by given any set, see if the sum of any collection of pairwise disjoint subsets of that set is greater than the set itself. This procedure is done for all sets to finally find a final solution that is the optimal one which give the highest revenue to the auctioneer.

Coalition structure generation is a complete analogy to combinatorial auction but from another perspective. It can be translated in such that the assets become agents, the sets of assets subjected to a bid is called a coalition and finally the

bid is the value the coalition is worth. Analogous to the previous example, the concept of, given N agents, generate a coalition such that the outcome is optimal in regards to a given property. This however will be an extreme challenge as the number of possible partitions grows exponential, specifically $O(n^n)$ [2]. One algorithm to solve this problem in $O(n^3)$ time is called the Dynamic Programming (DP) algorithm [3] described below.

Algorithm 1 Dynamic Programming algorithm

INPUT: $b(C)$ the bids for all sets $C \subseteq A$ where A is the set of assets.

VARIABLES: f a set that maps from a subset $C \subseteq A$ to a value

O a set that maps from a subset $C \subseteq A$ to the subset that maximizes the value for set C .

```

1: for all  $x \in A$ , do  $f(\{x\}) := b(\{x\})$ ,  $O\{x\} := \{x\}$  end for
2: for  $i := 2$  to  $n$  do
3:   for all  $C \subseteq A : |C| == i$  do
4:      $f(C) := \max\{f(C \setminus C') + f(C') : C' \subseteq C \wedge 1 \leq |C'| \leq \frac{|C|}{2}\}$ 
5:     if  $f(C) \geq b(C)$  then  $O(C) := C^*$  Where  $C^*$  maximizes right hand side of
       line 4 end if
6:     if  $f(C) < b(C)$  then  $f(C) := b(C) \wedge O(C) := C$  end if
7:   end for
8: end for

```

0.1.2 NVIDIA CUDA

CUDA (formerly known as Compute Unified Device Architecture) is a framework by NVIDIA equivalent to OpenCL which allows for C, C++ and FORTRAN code to be compiled and run on its Graphics Processing Units (GPUs). As GPUs are many-core processors designed to handle large amount of independent data when running a graphics program, it is advantageous to utilize such architecture to handle general purpose algorithms to solve complex problems with parallel nature. This architecture is described as Single Instruction Multiple Threads (SIMT) which could be seen similar to the Single instruction, multiple data (SIMD) architecture. The fundamental difference is that with SIMT, branch divergence is allowed.

The CUDA architecture consists of a scalable array of Streaming Multiprocessors (SM), where each SM contains from 8 cores to 192 programmable cores depending on the GPU version. Where one GPU can have several SMs. When launching a kernel which can be seen as a function call, the threads assigned are divided up into blocks of threads. Each block gets assigned to an SM to be pro-

cessed resulting in two blocks from the same kernel can be running concurrently on two different SMs. If there is no SM freely available the block will be put into a queue. The next thing to note is that each block is divided into warps of 32 threads which is scheduled by the warp scheduler. The reason is to hide memory read and write latency as while one warp is blocked waiting for data, another warp can run its instructions.

0.2 Literature review

With Combinatorial auctions and multiagent coalition formation being combinatorial optimization problem, an literature review was conducted to find relevant publications in the area. It was parted in two ways,

Using the GPU together with dynamic programming have been used before to solve similar combinatorial optimisation problems. Boyer, V et al successfully implemented and solved the knapsack problem with a great reduction in runtime speed [1]. Also introduced was ways of reducing memory occupancy which enable computation of much larger data sets as well as reducing the bandwidth utilized allowing for a faster processing speed.

0.3 Work

To this point I have implemented an algorithm called the Dynamic Programming (DP) algorithm by Rothkopf, M.H et al [3] which is used to solve combinatorial auctions, but it can also be used for complete coalition formation. It was first implemented as a sequential algorithm to allow firstly for a baseline to compare execution speed by. But also to gain knowledge about the algorithm to allow easier transition towards a parallel version. The algorithm I have derived that run on the GPU consists of generating a disjoint subset from a given set, and then calculated the value the disjoint sets give.

Input

f	The array which holds the bids
O	The step array that holds set configurations
Φ	A reference to an maximum value bucket unique for each kernel
C	The set to do subset construction on
Ψ	The maximum subset construction index

Variables

λ	How many bids should be evaluated per thread
Υ	A shared array containing threads maximum bid values
Δ	A shared array containing threads maximum bid subsets

v	A local array containing one of the threads bid value
C'	A local array containing one of the threads bid subsets
$bid = blockIdx.x$	Which block the threads belong to
$bdim = blockDim.x$	How many threads inside the block
$tid = threadIdx.x$	The thread index inside the block
$\psi := \lambda * (tid + bdim * bid)$	Initial subset construction index
Start of algorithm	
1: $\Delta_{tid} := \Upsilon_{tid} := 0$	Initilize values to zero
2: if $\psi \leq \Psi$ then	
3: for $i := 1$ to λ do	
4: $C'_i := (\neg C + \psi) \cap C$	Integer subset construction
5: if $ C'_i \leq \frac{ C }{2} \wedge \psi \leq \Psi$ then	
6: $v_i := f[C/C'_i] + f[C'_i]$	Set the split subset coalition value
7: end if	
8: if $v_i > \Upsilon_{tid}$ then	
9: $\Upsilon_{tid} := v_i$	Set value if v_{tid} is greatest so far
10: $\Delta_{tid} := C'_i$	Set subset if v_{tid} is greatest so far
11: end if	
12: $\psi := \psi + 1$	Increment by one
13: end for	
14: end if	
15: <i>syncthread</i> s	make sure all threads have finished
16: for $i := bdim >> 1$ to 0 do	
17: if $tid < i$ then	
18: if $\Upsilon_{tid} \leq \Upsilon_{tid+i}$ then	
19: $\Delta_{tid} := \Delta_{tid+i}$	
20: $\Upsilon_{tid} := \Upsilon_{tid+i}$	
21: end if	
22: end if	
23: <i>syncthread</i> s	make sure all threads have finished
24: $i := i >> 1$	bitshift i once to the right
25: end for	
26: if $tid == 0$ then	
27: if $\Upsilon_0 == 0$ then <i>exit</i> end if	exit if the bid is zero
28: if $f_C > \Upsilon_0$ then <i>exit</i> end if	exit if the bid is less than for C
29: if $atomicMax(\Phi, \Upsilon_0) < \Upsilon_0$ then	
30: $O_C = \Delta_0$	Sets the current optimal step in global memory
31: $f_C = \Upsilon_0$	Sets the current optimal value in global memory
32: <i>threadfence</i>	make sure all threads see the changes
33: end if	

```
34: end if  
35: return  $O, f$ 
```

Line 4 represent how the subset construction is made, it starts by taking the inverse of the integer C , add together the subset construction index ψ , and finally apply a bitwise AND with C to get a guaranteed subset of C . This implementation suffers though from generating subsets of greater cardinality needed, but testing have shown that the instruction overhead for generating these unwanted subsets only account for 18% to 8% overhead as shown by figure 2.

The code from 16 to 25 details how to do a reduction among all threads inside a block to find the maximum value from all of the subset constructions the threads have previously done. Where for each cycle of the loop, half of the threads will compare its value to a corresponding inactive thread and set its value and subset to whichever value is largest. Then it will reduce the number of active threads by half and continues until there is no more threads to be active. The loop will run $\log(n)$ times where n is the number of threads in the block.

Lines 26 to 34 represent how each block of a kernel update their maximum value and subset to the global memory for all the consecutive kernels to see. Beginning at line 26, here only the thread with an id of zero will pass in order to reduce the memory bandwidth but also as there is only need for one thread to update the global memory. Line 29 ensures that only a bigger value can be set by first querying the bucket for this kernel, if it is bigger it can continue to set the global memory. The reason for the atomic function is due to that blocks from the same kernel can and will run concurrently, therefore it will cause race conditions.

0.3.1 Result

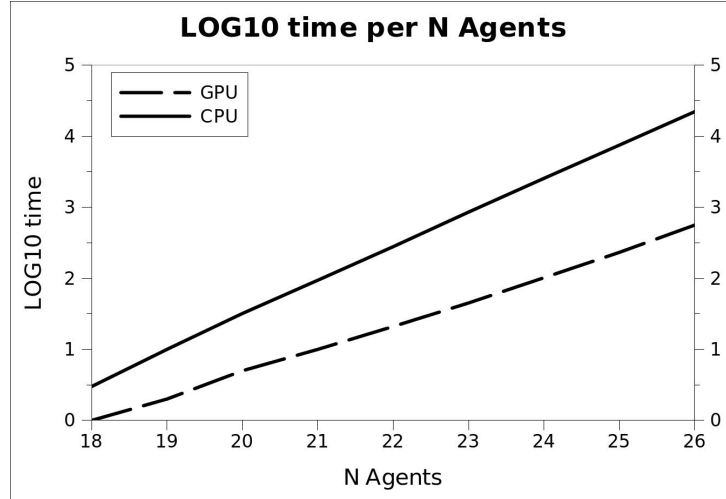


Figure 1: LOG10 runtime per N Agents

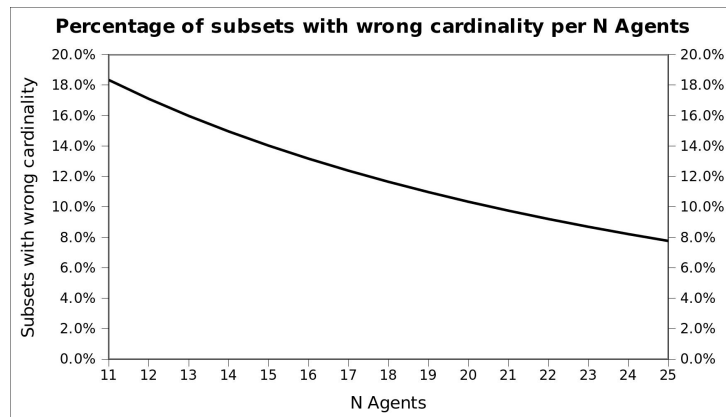


Figure 2: Percentage of generated subset of wrong cardinality

0.3.2 Evaluation

As it stands now, the algorithm that is running using the CUDA framework have an substantial performance increase as shown in figure 1. As seen from the figure, the slope of line CPU and GPU is different, indicating an different growth rate for an increase of n, thus having a different time complexity. Empirical measurement show that for each increase in n, the runtime for the cpu increase by three. While

the GPU bound algorithm starts its growth at 2 and the growth increases to 2.42 for 25 agents.

Figure 2 is the evaluation of line 4 on how many subsets generated that invalidates the condition $|C'_i| \leq \frac{|C|}{2}$ where $C'_i \subset C$. This is a measure that describes the instruction overhead for that particular part of the algorithm, and as it stands now with the algorithm being limited by memory latency due to random access, this instruction overhead can be considered negligible. This measure however does not take into account any duplicate subsets that are occasionally generated.

0.4 Future work

As there is currently a selected number of algorithms that is regarded faster than the DP algorithm such as the IDP algorithm [2] that have shown a promising efficiency. The next step would be to implement it both on the CPU and GPU. However as it is fairly similar to the DP algorithm, implementation may be trivial and further goals need to be determined. As the current work only works for a small number of agents, the next step would be to implement an algorithm to solve more realistic combinatorial auctions with a relative large number of agents and low number of coalition structures.

Bibliography

- [1] V. Boyer, D. El Baz, and M. Elkihel. Solving knapsack problems on gpu. *Computers & Operations Research*, 39(1):42–47, 2012.
- [2] T. Rahwan and N.R. Jennings. An improved dynamic programming algorithm for coalition structure generation. In *Proceedings of the 7th international joint conference on Autonomous agents and multiagent systems-Volume 3*, pages 1417–1420. International Foundation for Autonomous Agents and Multiagent Systems, 2008.
- [3] M.H. Rothkopf, A. Pekeč, and R.M. Harstad. Computationally manageable combinational auctions. *Management science*, 44(8):1131–1147, 1998.