

# Scheduling on Multi-Cores with GPU

Safia Kedad-Sidhoum<sup>1</sup>, Florence Monna<sup>1</sup>, Grégory Mounié<sup>2</sup>,  
Denis Trystram<sup>2,3</sup>

<sup>1</sup>Laboratoire d'Informatique de Paris 6, 4 Place Jussieu, 75005 Paris

<sup>2</sup>Grenoble Institute of Technology, 51 avenue Kuntzmann, 38330 Montbonnot  
Saint Martin, France

<sup>3</sup>Institut Universitaire de France

August 26, 2013

# Outline

# Scheduling with GPU

Most computers today include

- a Multi-core CPU
- a high performance parallel computing accelerators : the GPGPU (General Purpose Graphical Processing Unit).

Examples:

- Laptop/Tablet/Smartphone (Intel Core i7, Nvidia Tegra 4)
- Game console (PS4, Xbox One)
- Titan (in the top of the Top500 of the supercomputers)

In each machine, there are vectorial coprocessors with very high computing throughput, an interesting asset for High Performance Computing (HPC).

# GPU programming example

## Vector addition element by element

Compute  $Y = \alpha + X$ ,  $Y$  and  $X$  being two vectors of 1024 float.

```
1 prog = create_program([<<EOF
2 __kernel void addition( float alpha ,
3                          __global const float *x,
4                          __global float *y) {
5     size_t ig = get_global_id(0);
6     y[ig] = alpha + x[ig];
7 }
8 EOF
9 ])
10 create_kernel("addition", prog)
```

```
1 input = OpenCL::VArray::new(FLOAT, 1024)
2 output = OpenCL::VArray::new(FLOAT, 1024)
3 input_gpu = create_buffer(1024*4)
4 output_gpu = create_buffer(1024*4)
```

# GPU programming example

## Vector addition element by element

Compute  $Y = \alpha + X$ ,  $Y$  and  $X$  being two vectors of 1024 float.

```
1 prog = create_program([<<EOF
2 __kernel void addition( float alpha ,
3                         __global const float *x,
4                         __global float *y) {
5     size_t ig = get_global_id(0);
6     y[ig] = alpha + x[ig];
7 }
8 EOF
9 ])
10 create_kernel("addition", prog)
```

```
1 input = OpenCL::VArray::new(FLOAT, 1024)
2 output = OpenCL::VArray::new(FLOAT, 1024)
3 input_gpu = create_buffer(1024*4)
4 output_gpu = create_buffer(1024*4)
```

# GPU programming example

## Sequence of commands

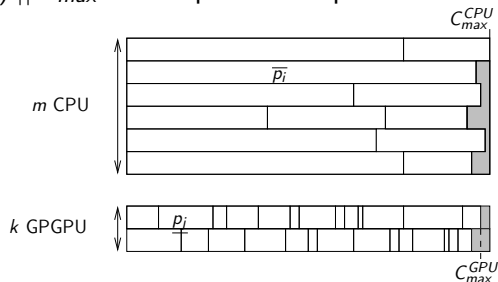
- line 1 copy input buffer from CPU memory to GPU memory
- line 2-4 compute the kernel with the arguments and vector of size 1024 float split in 64 bloc size
- line 5 copy output buffer from GPU memory to CPU memory

```
1 enqueue_write_buffer(1024*4, input, input_gpu)
2 args= set_args([OpenCL::Float::new(5.0),
3               input_gpu, output_gpu])
4 enqueue_NDrange_kernel(prog, args, [1024], [64])
5 enqueue_read_buffer(1024*4, output_gpu, output)
```

- The tasks assigned to the GPUs must be carefully chosen.
- Generic method to do the assignment for High Performance Computing Systems.
- No previous model:  
start with a simplified problem, without communication issues, precedence relations...

# Description of the Problem - Complexity

- $(Pm, Pk) \parallel C_{max}$ :  $n$  independent sequential tasks  $T_1, \dots, T_n$ .

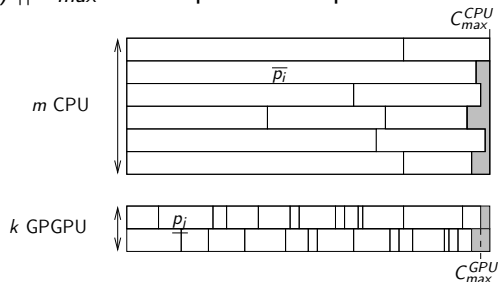


- Objective: minimize the **makespan** of the schedule.
- If  $\underline{p}_j = \bar{p}_j$  for all tasks  $(Pm, P1) \parallel C_{max} \Leftrightarrow P \parallel C_{max}$ , NP-hard  
 $\Rightarrow$  Problem of scheduling with GPUs also **NP-hard**



# Description of the Problem - Complexity

- $(Pm, Pk) \parallel C_{max}$ :  $n$  independent sequential tasks  $T_1, \dots, T_n$ .

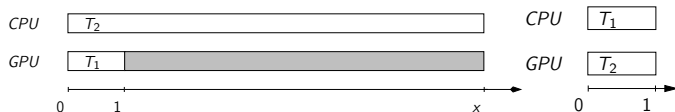


- Objective: minimize the **makespan** of the schedule.
- If  $\underline{p}_j = \bar{p}_j$  for all tasks  $(Pm, P1) \parallel C_{max} \Leftrightarrow P \parallel C_{max}$ , NP-hard  
 $\implies$  Problem of scheduling with GPUs also **NP-hard**

# List based scheduling

## Lemma

$(P1, P1) \parallel C_{max}$ : list scheduling algorithm has a ratio larger than the maximum speedup ratio of a task.



# Dual approximation technique

- Use of the dual approximation technique [Hochbaum & Shmoys, 1988]: for a ratio  $g$ , take a guess  $\lambda$  and either delivers a schedule of makespan at most  $g\lambda$ , or answers that there exists no schedule of length at most  $\lambda$ .
- At each step of the dual approximation, dynamic programming algorithm.
  - Case  $k = 1$ : performance ratio of  $g = \frac{4}{3}$  in time  $O(n^2 m^2)$ .
  - Case  $k \geq 2$  ratio  $g = \frac{4}{3} + \frac{1}{3k}$  in time  $O(n^2 m^2 k^3)$ .

# Dual approximation technique

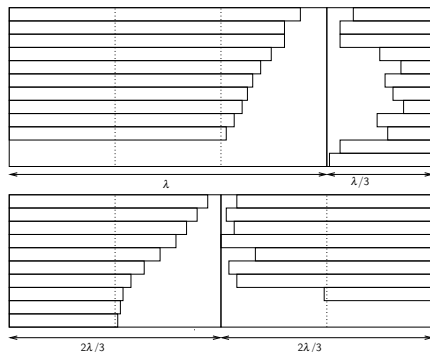
- Use of the dual approximation technique [Hochbaum & Shmoys, 1988]: for a ratio  $g$ , take a guess  $\lambda$  and either delivers a schedule of makespan at most  $g\lambda$ , or answers that there exists no schedule of length at most  $\lambda$ .
- At each step of the dual approximation, dynamic programming algorithm.
- Case  $k = 1$ : performance ratio of  $g = \frac{4}{3}$  in time  $O(n^2 m^2)$ .
- Case  $k \geq 2$  ratio  $g = \frac{4}{3} + \frac{1}{3k}$  in time  $O(n^2 m^2 k^3)$ .

# Dual approximation technique

- Use of the dual approximation technique [Hochbaum & Shmoys, 1988]: for a ratio  $g$ , take a guess  $\lambda$  and either delivers a schedule of makespan at most  $g\lambda$ , or answers that there exists no schedule of length at most  $\lambda$ .
- At each step of the dual approximation, dynamic programming algorithm.
- Case  $k = 1$ : performance ratio of  $g = \frac{4}{3}$  in time  $O(n^2 m^2)$ .
- Case  $k \geq 2$  ratio  $g = \frac{4}{3} + \frac{1}{3k}$  in time  $O(n^2 m^2 k^3)$ .

# The Shelves' Idea

For  $k = 1$ , assuming a schedule of length lower than  $\lambda$  exists.  
The idea: to partition the set of tasks on the CPUs into two sets, each consisting in two shelves.



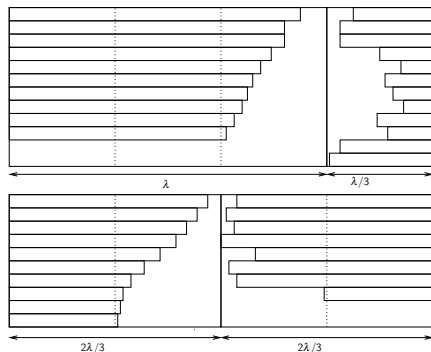
- a first set with a shelf of length  $\lambda$  the other of length  $\frac{\lambda}{3}$ ,
- a second set with two shelves of length  $\frac{2\lambda}{3}$ .

# The Shelves' Idea

- Partition ensures that the makespan on the GPUs is lower than  $\frac{4\lambda}{3}$ .
- The tasks are independent: the scheduling is straightforward when the assignment of the tasks has been determined.
- The main problem is to assign the tasks in each shelf on the CPUs or on the GPUs in order to obtain a feasible solution.

# Structure of an Optimal Schedule for $k = 1$

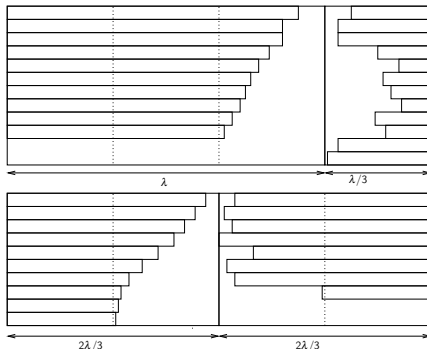
If there exists a schedule of length at most  $\lambda$





# Structure of an Optimal Schedule for $k = 1$

If there exists a schedule of length at most  $\lambda$

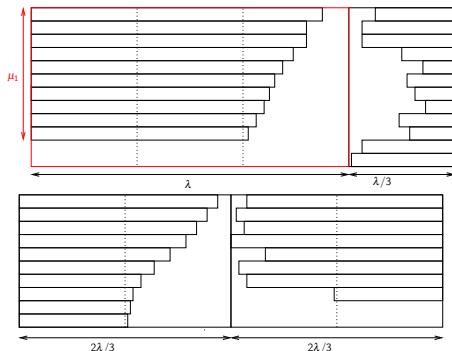


## Property (1)

For each task  $T_j$ ,  $\bar{p}_j \leq \lambda$ , and  $\sum_{\pi(j) \in \mathcal{C}} \bar{p}_j \leq m\lambda$ .

# Structure of an Optimal Schedule for $k = 1$

If there exists a schedule of length at most  $\lambda$

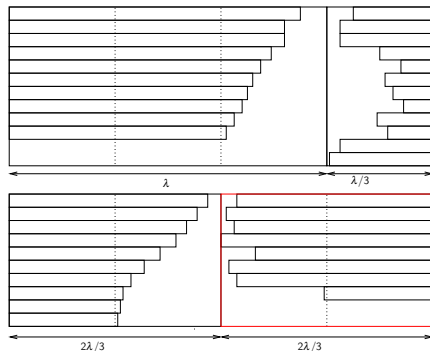


## Property (2)

$T_i, T_j$  two successive tasks on a CPU. If  $\bar{p}_i > \frac{2\lambda}{3}$ , then  $\bar{p}_j \leq \frac{\lambda}{3}$ .

# Structure of an Optimal Schedule for $k = 1$

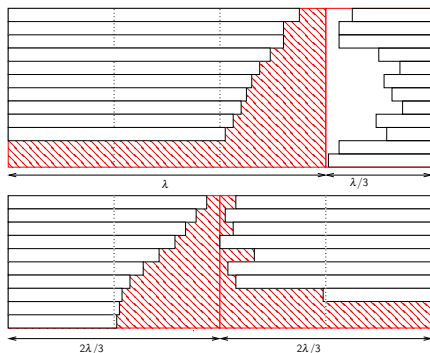
If there exists a schedule of length at most  $\lambda$



## Property (3)

Two tasks  $T_i, T_j$  with  $\frac{\lambda}{3} < \bar{p}_l \leq \frac{2\lambda}{3}$  ( $l = i, j$ ) can be executed successively on the same CPU within a time  $\frac{4\lambda}{3}$ .

# Structure of an Optimal Schedule for $k = 1$



The remaining tasks (with a processing time lower than  $\frac{\lambda}{2q}$  (resp.  $\frac{\lambda}{2q+1}$ )) fit in the remaining space in front of  $S_1$  and between all the others shelves, otherwise the schedule would not satisfy Property 1.

# Partitioning the Tasks into Shelves

We solve the assignment problem with a dynamic programming summing up the previous constraints. Here, we take  $g = \frac{4}{3}$ .

For task  $T_j$ , binary variable :  $x_j = \begin{cases} 1 & \text{if assigned to a CPU} \\ 0 & \text{if assigned to the GPU} \end{cases}$

$$W_C^* = \min \sum_{j=1}^n \bar{p}_j x_j \quad (1)$$

$$\text{s.t. } \frac{1}{2} \sum_{2\lambda/3 \geq \bar{p}_j > \lambda/3} x_j + \sum_{\bar{p}_j > 2\lambda/3} x_j \leq m \quad (2)$$

$$\sum_{j=1}^n \bar{p}_j (1 - x_j) \leq \frac{4\lambda}{3} \quad (3)$$

$$x_j \in \{0, 1\} \quad (4)$$

# Partitioning the Tasks into Shelves

- Dynamic programming algorithm: solves the previous problem in  $\mathcal{O}(n^2 m^2)$ .
- Reduction of the states on the GPU to a smaller number.  
Number of time intervals of length  $\frac{\lambda}{3n}$  for a task  $T_j$  executed on the GPU,  $v_j = \left\lfloor \frac{p_j}{\lambda/(3n)} \right\rfloor$ .  
 $N = \sum_{\pi(j) \in \mathcal{G}} v_j$  total number of these intervals on the GPU.
- Error on processing time of each task  $\varepsilon_j = \underline{p_j} - v_j \leq \frac{\lambda}{3n}$   
If all the tasks are assigned to the GPU, error at most  $n \frac{\lambda}{3n} = \frac{\lambda}{3}$ .
- Constraint (3) becomes  $N = \sum_{\pi(j) \in \mathcal{G}} v_j \leq 3n$

# Binary Search - Cost Analysis

- If optimum  $W_C^* = \min_{0 \leq \mu \leq m, 0 \leq \mu' \leq 2(m-\mu), 0 \leq N \leq 2n} W_C(n, \mu, \mu', N) > m\lambda$ , no solution with makespan  $\leq \lambda$  exists, algorithm answers "NO"
- Otherwise, construct **feasible solution** with makespan  $\leq \frac{4\lambda}{3}$ , with shelves on the CPUs and  $\mu^*$ ,  $\mu'^*$  and  $N^*$  values.
- One step of the dual-approximation algorithm, with a fixed guess. **Binary search** in  $\log(B_{max} - B_{min})$ .
- At each step, we have  $1 \leq j \leq n$ ,  $1 \leq \mu \leq m$ ,  $1 \leq \mu' \leq 2(m - \mu)$ , and  $0 \leq N \leq 3n$  so the time complexity is in  $O(n^2 m^2)$ .

Algorithm can be extended to  $(Pm, Pk) \parallel C_{max}$  with  $k \geq 2$ .

$$W_C^* = \min \sum_{j=1}^n \bar{p}_j x_j \quad (5)$$

$$\text{s.t. } \frac{1}{2} \sum_{2\lambda/3 \geq \bar{p}_j > \lambda/3} x_j + \sum_{\bar{p}_j > 2\lambda/3} x_j \leq m \quad (6)$$

$$\frac{1}{2} \sum_{2\lambda/3 \geq \underline{p}_j > \lambda/3} x_j + \sum_{\underline{p}_j > 2\lambda/3} x_j \leq k \quad (7)$$

$$N = \sum_{\pi(j) \in \mathcal{G}} v_j \leq 3kn \quad (8)$$

$$x_j \in \{0, 1\} \quad (9)$$



- The approximation algorithm can be extended to the problem with  $k \geq 2$  GPUs with a performance guarantee of  $\frac{4}{3} + \frac{1}{3k}$ .
- To solve each step of the binary search,  $O(n^2 k^3 m^2)$  states are considered, since  $1 \leq j \leq n$ ,  $1 \leq \mu \leq m$ ,  $1 \leq \mu' \leq 2(m - \mu)$ ,  $1 \leq \kappa \leq k$ ,  $1 \leq \kappa' \leq 2(k - \kappa)$ , and  $0 \leq N \leq 3kn$ .  
 $\implies$  Time complexity in  $O(n^2 k^3 m^2)$  for each step of the binary search.

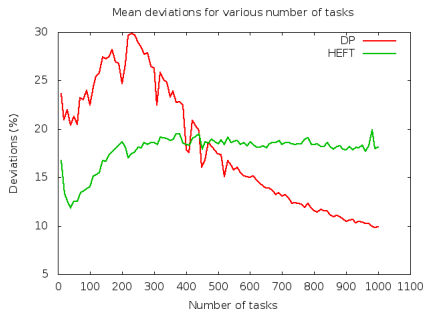
Comparison of a relaxed version of the dynamic programming (DP) algorithm with a ratio of 2 to the HEFT algorithm (Heterogeneous-Earliest-Finish-Time) [Topcuoglu et al. 2002]: prioritizing with decreasing average execution time, heterogeneous earliest finish time rule.

## Lemma

*For the  $(P_m, P_1)$ -problem, the worst case performance ratio of HEFT is larger than  $m/2$ .*

# Series of experiments

- Random instances: 10, 20, 40, 80 tasks, 1, 2, 4, 8, 16, 32, 64 CPUs, 1,2, 4, 8 GPUs.
- A task is assigned an acceleration factor of  $1/15$  or  $1/35$  with a probability of  $1/2$ .



# Conclusion and Perspectives

- Contribution: fast algorithms with constant guarantee for independent tasks on CPUs and GPUs. In the case of a single (resp. multiple) GPU(s) a ratio of  $\frac{4}{3} + \varepsilon$  (resp.  $\frac{4}{3} + \frac{1}{3k} + \varepsilon$ ) is achieved, and can be degraded to a ratio of 2 for efficiency.
- Finer ratios can be obtained with a sacrifice in time complexity.
- Extensions to partial preemption and malleable tasks can be done.
- On-going research on the problem with precedence relations, and data communication.
- Protocols in writing for an integration into parallel programming environment like StarPU and xKaapi.