

Scheduling on Hybrid Platforms: Improved Approximability Window

Vincent Fagnon¹ Imed Kacem²
Giorgio Lucarelli² **Bertrand Simon**³

1: LIG, Univ. Grenoble Alpes (France).

2: LCOMS, Univ. Lorraine (France).

3: IN2P3 Computing Center (France).

LATIN – January 2021

Scheduling on computing platforms

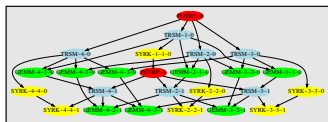
```
FOR k = 0..TILES-1
  FOR n = 0..k-1
    A[k][k] <- DSYRK(A[k][n], A[k][k])
    A[k][k] <- DPOTRF(A[k][k])
  FOR m = k+1..TILES-1
    FOR n = 0..k-1
      A[m][k] <- DGEWM(A[k][n], A[m][n], A[m][k])
    A[m][k] <- DTRSM(A[k][k], A[m][k])
```



Runtime Software

- ▶ Schedule tasks
- ▶ Aware of durations
- ▶ Objective: complete the graph ASAP (minimize the makespan)

Scheduling on computing platforms



Runtime Software

- ▶ Schedule tasks
- ▶ Aware of durations
- ▶ Objective: complete the graph ASAP (minimize the makespan)

Scheduling on computing platforms



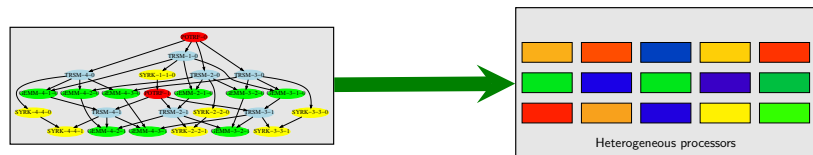
Runtime Software

- ▶ Schedule tasks
- ▶ Aware of durations
- ▶ Objective: complete the graph ASAP (minimize the makespan)

Classical models

- ▶ m identical processors: $P|prec|C_{max}$ - ignore accelerators (GPUs...)

Scheduling on computing platforms



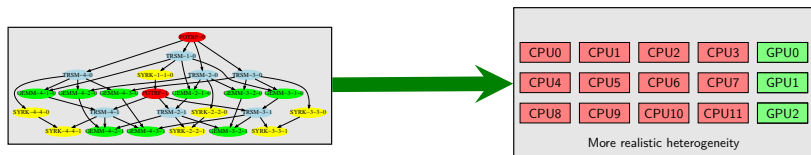
Runtime Software

- ▶ Schedule tasks
- ▶ Aware of durations
- ▶ Objective: complete the graph ASAP (minimize the makespan)

Classical models

- ▶ m identical processors: $P | prec | C_{max}$ - ignore accelerators (GPUs...)
- ▶ m unrelated processors: $R | prec | C_{max}$ - too complex

Scheduling on computing platforms



Runtime Software

- ▶ Schedule tasks
- ▶ Aware of durations
- ▶ Objective: complete the graph ASAP (minimize the makespan)

Classical models

- ▶ m identical processors: $P | prec | C_{max}$ - ignore accelerators (GPUs...)
- ▶ m unrelated processors: $R | prec | C_{max}$ - too complex

Our model

- ▶ Two types of processors: e.g., m CPUs and k GPUs

Input

- ▶ Graph of tasks with precedence constraints
- ▶ m identical processors (CPU) and k identical processors (GPU)
- ▶ $m \geq k$
- ▶ For each task: its running time \bar{p}_i on CPU and \underline{p}_i on GPU

Output

- ▶ Schedule of minimum makespan (ignoring communication times)

Evaluation metric

- ▶ Approximation factor: maximum value of $\frac{\text{obtained makespan}}{\text{optimal makespan}}$

Note: other scenarios studied (independent tasks, online arrival of tasks)

Two sub-problems

- ▶ **allocation**: on which processor type (CPU/GPU) we place each task
- ▶ **schedule**: once the allocation is fixed, when each task is run

Assume the allocation is fixed

Best known algorithm: **list-scheduling** [Graham '66]

(if a processor is available and a task can run on it, do it)

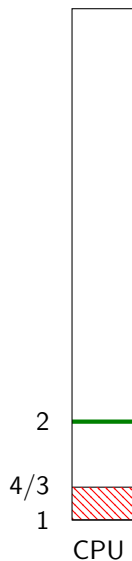
We can define:

- ▶ W_C (resp. W_G): total load on CPUs (resp. GPUs)
- ▶ CP : time to complete the longest path (critical path)

Lemma

List-scheduling makespan is at most $CP + \frac{W_C}{m} + \frac{W_G}{k}$ so it is a 3-approximation when the allocation is fixed.

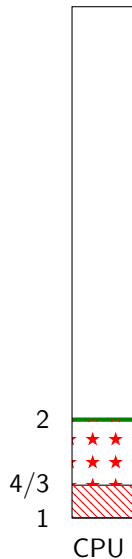
Approximation window – a) identical processors



List-Scheduling [Graham '66]

NP-hard [Lenstra & Rinnooy Kan '78]

Approximation window – a) identical processors

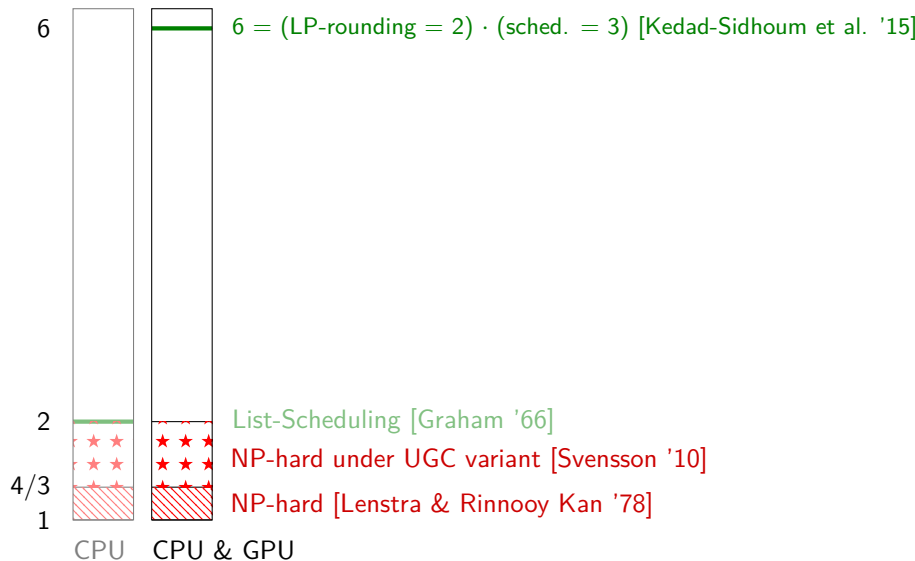


List-Scheduling [Graham '66]

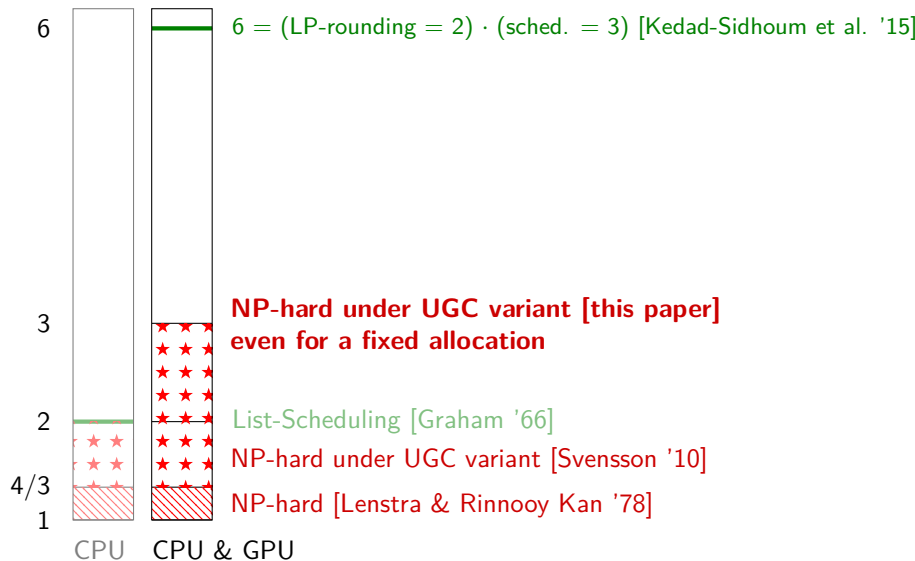
NP-hard under UGC variant [Svensson '10]

NP-hard [Lenstra & Rinnooy Kan '78]

Approximation window – b) CPUs & GPUs



Approximation window – b) CPUs & GPUs



Approximation window – b) CPUs & GPUs

6
 $3 + 2\sqrt{2}$

$6 = (\text{LP-rounding} = 2) \cdot (\text{sched.} = 3)$ [Kedad-Sidhoum et al. '15]

LP-rounding + subtile analysis [this talk]

3

**NP-hard under UGC variant [this paper]
even for a fixed allocation**

2

List-Scheduling [Graham '66]

$4/3$

NP-hard under UGC variant [Svensson '10]

1

NP-hard [Lenstra & Rinnooy Kan '78]

CPU

CPU & GPU

Approximation window – b) CPUs & GPUs

6
 $3 + 2\sqrt{2}$

$6 = (\text{LP-rounding} = 2) \cdot (\text{sched.} = 3)$ [Kedad-Sidhoum et al. '15]

LP-rounding + subtile analysis [this talk]

Open Question:

how much do we lose to decide the allocation?

3

**NP-hard under UGC variant [this paper]
even for a fixed allocation**

2

List-Scheduling [Graham '66]

NP-hard under UGC variant [Svensson '10]

$4/3$

NP-hard [Lenstra & Rinnooy Kan '78]

1

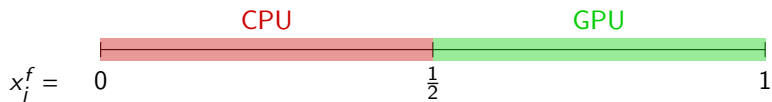
CPU

CPU & GPU

Idea: find an allocation by rounding an LP solution,
then use List Scheduling.

$$\begin{array}{l}
 \text{minimize: } C \\
 \frac{1}{m} \mathbf{W}_C \leq C \\
 \frac{1}{k} \mathbf{W}_G \leq C \\
 \mathbf{CP} \leq C
 \end{array}
 =
 \begin{array}{l}
 \text{minimize: } C \\
 \frac{1}{m} \sum_j \bar{p}_j x_j \leq C \\
 \frac{1}{k} \sum_j \underline{p}_j (1 - x_j) \leq C \\
 C_i + \bar{p}_j x_j + \underline{p}_j (1 - x_j) \leq C_j \text{ for all } i \rightarrow j \\
 0 \leq C_j \leq C \\
 x_j \in [0, 1]
 \end{array}$$

x_j : equals 0 (resp. 1) if task j goes to CPU (resp. GPU)



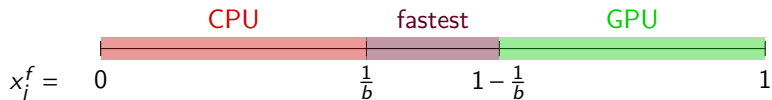
Hence, we have ($\cdot^f \rightarrow$ optimal fractional solution):

$$\begin{aligned}
 C_{max} &\leq \frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G + \mathbf{C}P \\
 &\leq 2 \cdot \frac{1}{m} W_C^f + 2 \cdot \frac{1}{k} W_G^f + 2 \cdot CP^f \\
 &\leq 6 \cdot OPT
 \end{aligned}$$

Tight in two ways:

- ▶ approximation factor = 6 (reached on an example)
- ▶ LP integrality gap = 2

Better rounding of the LP ($b > 2$)



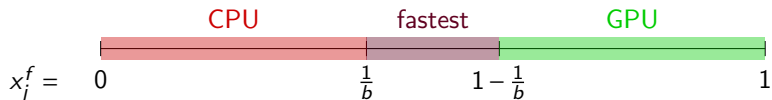
We gain on the critical path: $\mathbf{CP} \leq \frac{b}{b-1} \mathbf{CP}^f$,

and lose on the loads: $\frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G \leq b \cdot \left(\frac{1}{m} \mathbf{W}_C^f + \frac{1}{k} \mathbf{W}_G^f \right)$.

Hence,

$$\begin{aligned} C_{max} &\leq \frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G + \mathbf{CP} \\ &\leq \end{aligned}$$

Better rounding of the LP ($b > 2$)



We gain on the critical path: $\mathbf{CP} \leq \frac{b}{b-1} \mathbf{CP}^f$,

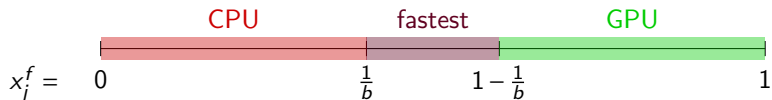
and lose on the loads: $\frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G \leq b \cdot \left(\frac{1}{m} W_C^f + \frac{1}{k} W_G^f \right)$.

But we have $\mathbf{W}_C + \mathbf{W}_G \leq \frac{b}{b-1} (W_C^f + W_G^f)$

Hence,

$$\begin{aligned} C_{max} &\leq \frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G + \mathbf{CP} \\ &\leq \end{aligned}$$

Better rounding of the LP ($b > 2$)



We gain on the critical path: $\mathbf{CP} \leq \frac{b}{b-1} \mathbf{CP}^f$,

and lose on the loads: $\frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G \leq b \cdot \left(\frac{1}{m} \mathbf{W}_C^f + \frac{1}{k} \mathbf{W}_G^f \right)$.

But we have $\mathbf{W}_C + \mathbf{W}_G \leq \frac{b}{b-1} (\mathbf{W}_C^f + \mathbf{W}_G^f)$

Hence,

$$\begin{aligned} C_{max} &\leq \frac{1}{m} \mathbf{W}_C + \frac{1}{k} \mathbf{W}_G + \mathbf{CP} \\ &\leq \frac{1}{m} (\mathbf{W}_C + \mathbf{W}_G) + \frac{m-k}{mk} \mathbf{W}_G + \mathbf{CP} \\ &\leq \dots \leq (3 + 2\sqrt{2}) \cdot OPT \text{ for } b = 1 + \sqrt{2} \end{aligned}$$

Same assumption as in [Bazzi, Norouzi-Fard '15], a variant of the UGC stronger than the one used in [Svensson '10, Bansal & Khot '09]

- ▶ Stays valid when the allocation is fixed
- ▶ Stays valid if processors are *related*: $\frac{\bar{p}_i}{p_i}$ is the same for all tasks
- ▶ Stays valid for unrelated processors and any value of m/k

Note: m/k linked to the online problem difficulty (best = $\Theta\left(\sqrt{\frac{m}{k}}\right)$)

[Amaris, Lucarelli, Mommessin, Trystram '19
Canon, Marchal, Simon, Vivien '19]

Approximation ratio in function of m/k

