

# Sort clusters algorithms

September 20, 2020

September 20, 2020

## 1 Introduction

Something about CAL language and a whole BaltiLSC system and the project basic information ...

## 2 Data structures

### 2.1 Cluster-based structures

```
1 class Machine:
2     MEMORY: float # [GB]
3     CPUs: List[float] # [GHz]
4     GPUs: List[float]
5     PRICE: float # [EURO/h]
6
```

Listing 1: Machine data structure

```
1 class Resources:
2     CPU: List[float] # [GHz]
3     GPU: List[float]
4     MEMORY: float # [GB]
5
```

Lets indroduce the structure that represents the class of a machine. It takes into account the generalized resources of the machine across their scope:

```
7 class MachineClass:
8     MIN: Resources
9     MAX: Resources
10
```

Listing 2: Machine data structure

```
1 class Storage:
2     SIZE_RANGE: Tuple[float, float] # [GB]
3     TYPE: Enum["SSD", "HDD"]
```

```

4     PRICE: float # [EURO/(TB*h)]
5     PricePerHour: float
6

```

```

11    class Cluster:
12        ID: int
13        MACHINES: Dict[str: Machine]
14

```

Listing 3: Cluster data structure

## 2.2 Application-based structures

```

1    class Data:
2        TYPE: Enum['SET_OF_IMAGES', 'CSV_FILE', ...]
3

```

```

1    class DataSpec(Data):
2        OVERALL_SIZE: float # [GB]
3        QUANTITY: int
4        LARGEST_PART_SIZE: float # [GB]
5

```

```

1    class Module:
2        NAME: str
3        REQUIRED_RESOURCES: Resources
4        INC: float # internal network capacity [Gb/s]
5        ENC: float # external network capacity [Gb/s]
6

```

Listing 4: Module data structure

```

1    class Application:
2        MODULES: Dict[str: Module]
3        INPUT_DATA: Data
4        OUTPUT_DATA: Data
5        CAL_DIAGRAM: str # CAL diagram in .yaml format
6        AIFR: func(data_size: float): float # Function that returns
        internal flow ration of the application based on input data
        size*
7

```

\*internal flow ratio will be describe soon

```

1    class Task:
2        APP: Application
3        INPUT_DATA: DataSpec
4        RESOURCES_RESERVATION_RANGE: Tuple[Resources, Rersources]
5

```

```

1    class TaskExecution:
2        TASK_ID: int
3        CLUSTER_ID: int
4        TIME: float # [s]
5

```

Listing 5: TaskExecution data structure

```

1 class ModuleExecution:
2     MODULE_ID: int
3     MACHINE_ID: int
4     TIME: float # [s]
5

```

Listing 6: ModuleExecution data structure

## 3 Algorithm implemetation

### 3.1 Input and output data

As a input data for algorithm

1. *task\_1*: Task - the task we want to execute,
2. *matching\_clusters*: List[Cluster] - list of clusters that meet resources reservation range of the *task\_1*
3. *time\_price\_ratio*: float [0,1] - the priority of sorting policy set up by user who is executing the task

As a output data of the algorithm we receive a list of sorted clusters with paired modules for individual machines classes. The example output can look like this (MachineClass, Cluster and Module simplified data structures are shown in listening 2, 3 and 4):

```

1 [
2     (
3         c_1: Cluster,
4         [
5             (m_1: Module, MC_1: MachineClass),
6             (m_2: Module, MC_1: MachineClass),
7             (m_3: Module, MC_4: MachineClass),
8             ...
9         ]
10    ),
11    (
12        ...
13    ),
14    ...
15 ]
16

```

Listing 7: Example output of the sorting algorithm

### 3.2 Algorithms utils

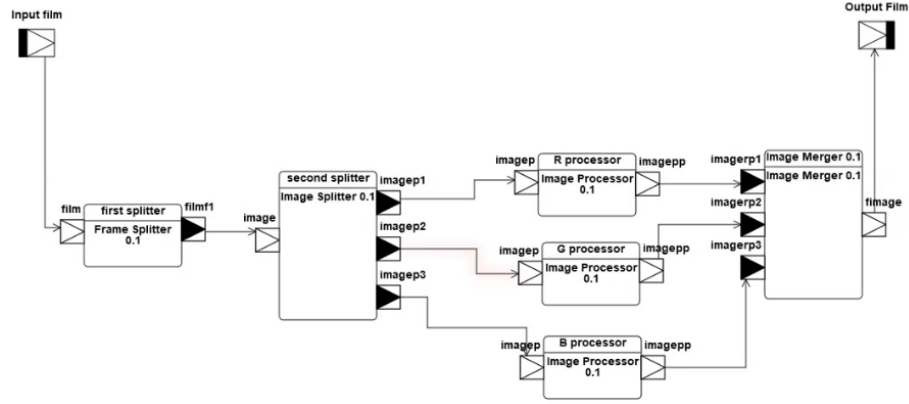
Let's get through our algorithms with an app as follows:

```

1 app_1 = Application(
2     MODULES = {
3         'first splitter': m1,
4         'second splittter': m2,

```

Figure 1: *app\_1* implementation with CAL language



```

5     'R processor': m3,
6     'G processor': m3,
7     'B processor': m3,
8     'Image Merger 0.1': m4,
9 },
10 INPUT_DATA = Data(
11     TYPE = 'FILM'
12 ),
13 OUTPUT_DATA = Data(
14     TYPE = 'FILM'
15 ),
16 CAL_DIAGRAM = '(call diagram content based on application
17 implementation)',
18 AIFR = func(int dataSize) { ... } # See the figure below
19 )

```

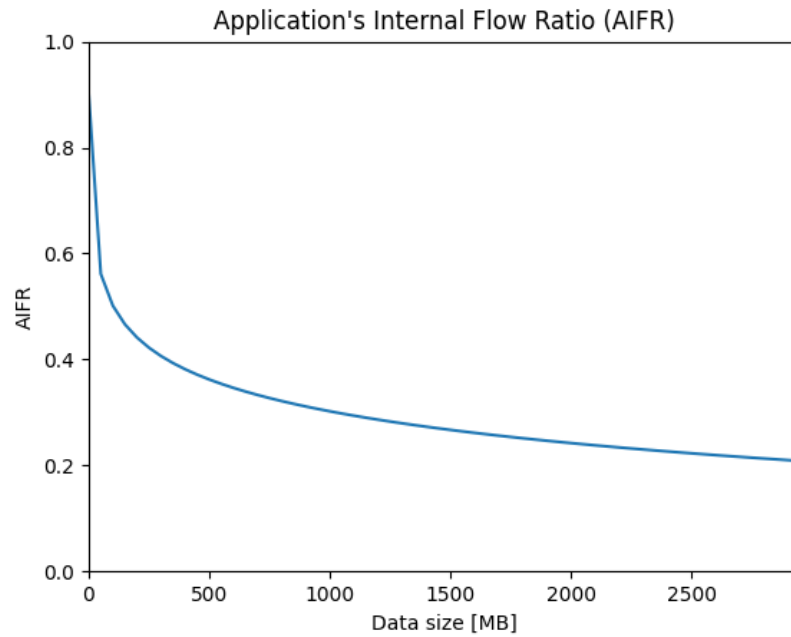
An example Task structure is listen below:

```

1 task_1 = Task(
2     APP = app_1,
3     INPUT_DATA: DataSpec(app_1.INPUT_DATA) (
4         OVERALL_SIZE = 1.2
5         QUANTITY = 1
6         LARGEST_PART_SIZE = 1.2
7     ),
8     RESOURCES_RESERVATION_RANGE: (
9         Resources(
10             CPU = [1.8],
11             GPU = [],
12             MEMORY = [2],
13         ),
14         Resources(
15             CPU = [2.4, 2.4, 2.4, 2.4, 2.4, 2.4],
16             GPU = [],
17             MEMORY = [4],
18         )
19 )

```

Figure 2: *app\_1 AIFR*



20 )  
21

Let's assume that we have for a disposition the following clusters:

```
1 c1: Cluster = Cluster(  
2   ID: 1,  
3   MACHINES: {  
4     '1': Machine = (  
5       MEMORY = 2  
6       CPUs = [1.8]  
7       GPUs = []  
8       PRICE = [0.12]  
9     ),  
10    '2': Machine = (  
11      MEMORY = 4  
12      CPUs = [2.2]  
13      GPUs = []  
14      PRICE = [0.18]  
15    )  
16  },  
17 )  
18 c2: Cluster = Cluster(  
19   ID: 2,  
20   MACHINES: {  
21     '1': Machine = (  
22       MEMORY = 2
```

```

23     CPUs = [1.8]
24     GPUs = []
25     PRICE = [0.15]
26 ),
27 '2': Machine = (
28     MEMORY = 4
29     CPUs = [2.4]
30     GPUs = []
31     PRICE = [0.25]
32 )
33 },
34 )
35

```

The helper functions used by sorting algorithms are listed below:

1. `getAppModulesLine(app: Application) → List[Module]` - return the sequence of app's modules

```

1     >>> getAppModulesLine(app_1)
2     [m1, m2, m3, m4]
3

```

### 3.3 Price-based algorithm

Price-based algorithm takes following assumptions:

1. Each instance of the same *MachineClass* is equal in terms of computation speed, regardless of the cluster it belongs to.
2. Internal data flow time is omitted.
3. The output is a list of sorted clusters without any information about computation times.
4. Algorithm estimates the relative (in relation to other clusters) costs of task execution and do not consider the computation times at all in the final result.

The minimum data set that we need to execute the algorithm consist of the specification for each machine (more specific: a machine resources to assign the appropriate class and the price per hour for using it).

Price-based algorithm steps:

1. `getModulesToMachineClassProbability(app: Application, clusters: List[Cluster], rrr: Tuple[Resources, Resources]) → Dict[int: List[Tuple[Module, List[Tuple[MachineClass, float]]]]]` - return the list of module to machine class pairs (with probabilities of each matching) for each cluster. It takes into consideration the resource reservation range and the resource requirements for each module. The probability can depends on how frequently the modules are assigned to each *MachineClass* (the extended version of the algorithm using historical data). In this simple example we assume uniform distribution of probabilities and that each module can be executed on any class of a machine.

```

1  >>> getModulesToMachineClassProbability(app_1, [c_1, c_2],
2      (resources_min, resources_max))
3  {
4      1: [
5          (m1, [
6              (mc_1, 0.5),
7              (mc_2, 0.5),
8          ]),
9          ...
10         (m4, [
11             (mc_1, 0.5),
12             (mc_2, 0.5),
13         ]),
14     ],
15     2: [
16         (m1, [
17             (mc_1, 0.5),
18             (mc_2, 0.5),
19         ]),
20         ...
21         (m4, [
22             (mc_1, 0.5),
23             (mc_2, 0.5),
24         ]),
25     ],
26 ]
27 }
28
29
30

```

2. In the list of modules to machine class probabilities for each cluster we change each machine class instance to its price per hour and receive:

```

1  {
2      1: [
3          (m1, [
4              (0.12, 0.5),
5              (0.18, 0.5),
6          ]),
7          ...
8          (m4, [
9              (0.12, 0.5),
10             (0.18, 0.5),
11         ]),
12     ],
13     2: [
14         (m1, [
15             (0.15, 0.5),
16             (0.25, 0.5),
17         ]),
18         ...
19         (m4, [
20
21
22

```

```

23         (0.15, 0.5),
24         (0.25, 0.5),
25     ]
26 ],
27 }
28

```

3. Next, we make internal multiplication of elements of List for each module:

```

1  {
2      1: [
3          (m1, 0.15),
4          ...
5          (m4, 0.15),
6      ]
7      2: [
8          (m1, 0.2),
9          ...
10         (m4, 0.2),
11     ]
12 }
13

```

4. Finally, we can sum the prices for each module execution with appropriate weights. Weights are calculated based on time execution ratios between modules and the size of the input data (in this simple example we assume equal execution times among all of the modules):

```

1  >>> getModulesTimeExecutionRatios(task_1: Task) -> List[
2      Tuple[Module, ratio]]
3      [(m1, 0.25), (m2, 0.25), (m3, 0.25), (m4, 0.25)]

```

The result of price-based algorithm based on our example:

```

1  [(1, 0.15), (2, 0.2)]
2

```

The algorithm accuracy could be increased by better methods of estimating time execution ratios. This methods will be described in the following parts.

### 3.4 Test-execution-based algorithm

No possession for the historical data is required by test-execution-based algorithm. It only uses data from the set of testing tasks executions. Lets say that we have 3 clusters we want to sort:

```

1  clusters_list: List[Cluster] = [c1, c2, c3]
2

```

Listing 8: List of available clusters

On the each of the clusters we run the same set of testing tasks consist of a single JobBatch.



```

1  testing_tasks: List[Task] = [t1, t2, t3]
2

```

Listing 9: Set of testing tasks

We end up with the testing tasks performances summaries:

```

1  testing_tasks_performances: List[Tuple(Cluster, Task, float,
2    float)] = [
3    (c1, t1, 0.21, 0.11),
4    (c1, t2, 0.44, 0.31),
5    (c1, t3, 0.24, 0.13),
6    (c2, t1, 0.66, 0.07),
7    (c2, t2, 0.84, 0.04),
8    (c2, t3, 0.54, 0.09),
9    (c3, t1, 0.39, 0.21),
10   (c3, t2, 0.57, 0.24),
11   (c3, t3, 0.33, 0.12),
12 ]

```

Listing 10: Testing tasks performances summaries

The two ending floats are respectively the execution time (in hours) and the execution price (in EUR). Having such an information we can count price and time clusters coefficients. For each testing task we search the minimal execution time and execution price. Then for each cluster we divide each time and price by the minimal values corresponding to appropriate tasks. We end up with the following coefficients:

```

1  testing_tasks_coefficients: List[Tuple(Cluster, Task, float,
2    float)] = [
3    (c1, t1, 1.0, 1.57),
4    (c1, t2, 1.0, 7.75),
5    (c1, t3, 1.26, 1.00),
6    (c2, t1, 3.14, 1.0),
7    (c2, t2, 1.91, 1.0),
8    (c2, t3, 2.84, 1.44),
9    (c3, t1, 1.86, 3.0),
10   (c3, t2, 1.3, 6.0),
11   (c3, t3, 1.0, 1.33),
12 ]

```

Listing 11: Testing tasks performances summaries

Next, we group the table above by cluster with sum as a aggregation function for the price and time coefficients.

```

1  cluster_coefficients: List[Tuple(Cluster float, float)] = [
2    (c1, 3.26, 10.32),
3    (c2, 7.89, 3.44),
4    (c3, 4.16, 10.33),
5  ]
6

```

Listing 12: Cluster coefficients

Finally, using the user typed time-price priority (a coefficient in range  $[0, 1]$ , let's assume it's 0.5 in our simple example) we can count the result:

```
1  [
2    (c1, 3.26*0.5 + 10.32*0.5),
3    (c2, 7.89*0.5 + 3.44*0.5),
4    (c3, 4.16*0.5 + 10.33*0.5),
5  ]
6  final_result = [
7    (c2, 5.67),
8    (c1, 6.79),
9    (c3, 7.25),
10 ]
11
```

Listing 13: any

### 3.5 Advance algorithm

The algorithm described above is based on execution of testing tasks and the historical data of task executions. To make it work we need to establish the set of well diversified testing applications and modules. The sorting is based on time ratios for the execution of test tasks, so the algorithms used by modules from test tasks should have linear computational complexity to make sorting most reliable.

Every new cluster wanting to join the fun should run each application with a few different input data sets. This data sets should be properly establish as well. Every new machine appering as available on the cluster should perform computations of all of the modules from the testing task set (if the machine conforms resources requirements of modules of course). Futhermore, we introduce the concept of machine class to generalize the calculation of machine performance.

After performing tests described above we will end up with testing executions data required for sorting algorithm implementation. In listenings 5 and 6 you can find simplified format of the testing tasks executions data.

Sorting has the following features:

1. Each cluster is tested for it internal data time flow beetwen nodes and the sorting is using this times,
2. Each machine on cluster is tested for it performance time and the sorting is using this times,
3. User before starting the task execution can provide time\_price\_ratio setting a priority on a more time or financial basis,
4. The algorithm not only indicate the most favorable cluster, but also pair the task modules to the appropriate machine classes.

Individual steps of the sorting algorithm:

1. `getTimeLines(task_1: Task) → List[List[Tuple[Module, int]]]` - this will return a list of sequence of modules that will be carried out in succession with the number representing the parallelism of the module. It will return the list of all possible time lines based on the configuration of the application and the `task_1.RESOURCES_RESERVATION_RANGE`.

```

1     >>> time_lines: List[Module] = getTimeLines(task_1)
2     >>> print(time_line)
3     [
4         [(m1, 1), (m2, 1), (m3, 1), (m4, 1)],
5         ...,
6         [(m1, 1), (m2, 4), (m3, 12), (m4, 1)]
7     ]
8

```

2. `getDetailedTimeLines(time_lines: List[List[Tuple[Module, int]]], task_1: Task, matching_clusters: List[Cluster]) → Dict[Tuple[Cluster, str]: List[Tuple[Module, int, MachineClass]]]` - return collection of modules list paired with machines (based on `matching_clusters` list and `task_1.RESOURCES_RESERVATION_RANGE`)

```

1     >>> detailed_time_lines: Dict[str: List[Tuple[Module,
2         Machine]]] = getDetailedTimeLines(time_line, task_1, [c1,
3         c2])
4     >>> print(detailed_time_line)
5     {
6         (c1, 'dl_0'): [(m1, 1, m_c_1), (m2, 1, m_c_1), (m3, 1,
7             m_c_1), (m4, 1, m_c_1)],
8         ...,
9         (c2, 'dl_0'): [(m1, 1, m_c_1), (m2, c2.MACHINES['1']), (
10             m3, 1, m_c_1), (m4, 1, m_c_1)],
11         ...,
12         (c2, 'dl_15'): [(m1, 1, m_c_2), (m2, 1, m_c_2), (m3, 1,
13             m_c_2), (m4, 1, m_c_2)],
14         ...,
15         (c2, 'dl_767'): [(m1, 1, m_c_2), (m2, 4, m_c_2), (m3,
16             12, m_c_2), (m4, 1, m_c_2)],
17     }
18

```

3. `getInternalFlowCoefficient(task_execution_1: TaskExecution, c1: Cluster) → float` - returns a time coefficient of internal data flow between modules for a given TaskExecution performed on cluster. To be more specific, it returns overall time of task execution minus sum of the maximum executions times of each modules from the task time line. This function using following formula:

$$coef = \frac{internal\_flow\_time}{task\_execution\_time} \quad (1)$$

In our example (*app\_1*, modules parallization respectively equals 1, 3, 4, 1) shown in Figure 3 the coefficient will eqaul:

$$coef = \frac{T - t1 - t2 - t3 - t4}{T} \quad (2)$$

So as a result we get the internal flow coefficient ( $IFC$ ) in range  $[0, 1]$  the

Figure 3: *app\_1* timeline

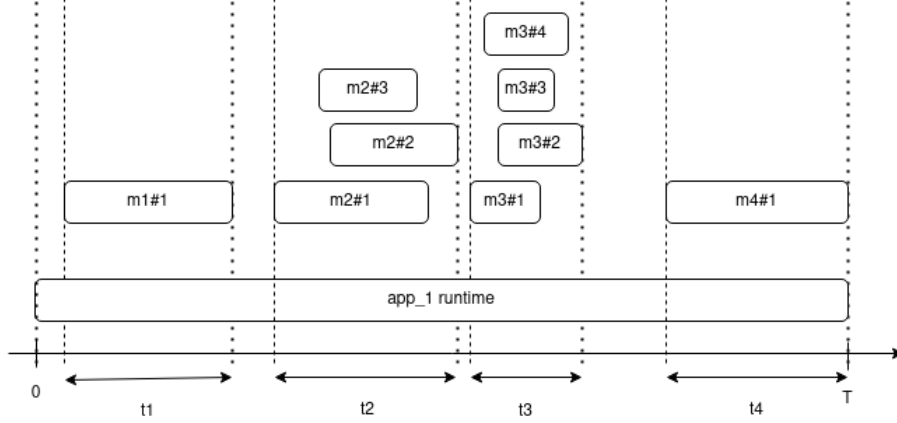
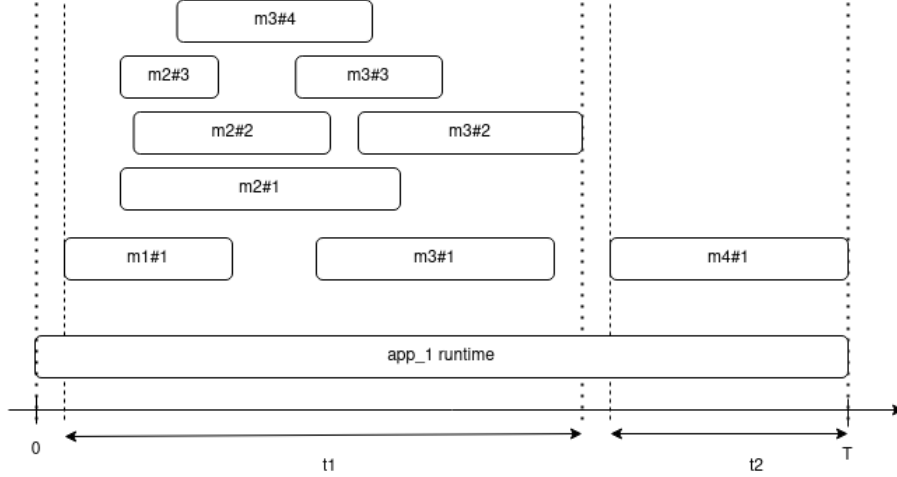


Figure 4: *app\_1* timeline with parallelization beetwen modules



value of which reflects the ratio of internal data flow between modules time and the task execution time. If there is no need for modules to be executed in sequence (one after another), it is possible that later modules will be executed in parallel with the earlier ones. This is shown in Figure 4. In such a case the  $IFC$  will be significantly lower. Even though, it is possible to count the  $IFC$  for such a cases, the algorithm imply the modules are executed one after another.

4. `getInternalFlowsRatios(test_tasks_executions: SetTaskExecution, task_1:`

Task, clusters: Set[Cluster] → Dict[Cluster: float] - count the average internalFlowCoefficient for each cluster from clusters set and returns the collection of ratios of received coefficients. The final ratios are the result of multiplying normalized internal flow coefficients by the task internal flow coefficient.

```

1     >>> internal_flows_ratios = getInternalFlowsRatios(
2     test_tasks_executions, [c1, c2])
3     >>> print(internal_flows_ratios)
4     {
5         c1: 1.19,
6         c2: 1.0
7     }

```

5. getTestingModulesExecutionTimesRatios(test\_tasks\_executions: List[TaskExecution], details\_time\_lines: Dict[Tuple[Cluster, str]: List[Tuple[Module, int, MachineClass]]) → Dict[str: List[float]] - we take first Tuple[Module, int, MachineClass] from each element of details\_time\_lines. Then we create an unique set of machine classes with number of parallization divided (with rounding up) by available parallization based on machine class resources range (so if we have 4 available CPUs for a machine class and the parallization coefficient for the modul under consideration is 10, then have to multiply module execution time by roundUp(10/4) = 3). In our example it will be just:

```

1     machines_1 = { (m_c_1, 1), ..., (m_c_1, 1), ..., (m_c_1,
2     1), ..., (m_c_1, 1) }

```

So for the first Tuple[Module, int, MachineClass] the rations will equal 1. Based on test\_tasks\_executions we can extract modules executions common to set of machines described above. Lets say it will be a following set:

```

1     m_1_mod_execs: Set[TaskExecution] = {t1m1, ..., t1m4, ...,
2     t5m2, ..., t6m1}

```

Having information about the execution times of modules on machine classes, for each module we count execution time ratio for each machine class (taking into consideration the number of parallization for each module):

$$t1m1.c1m1 = \frac{t1m1.TIME / \#\{module\_parallization\}}{[t1m.TIME \text{ for } t1m \text{ in } m\_1\_mod\_execs].sum()} \quad (3)$$

So for each machine class we will receive a set of modules executions times ratios:

```

1     c1m1: Set[float] = {t1m1_c1_m_c_1, ..., t1m4_c1_m_c_1,
2     ..., t5m2_c1_m_c_1, ..., t6m1_c1_m_c_1}

```

Then we count the average ratio of testing modules executions for each class machine:

```
1   c1m1_ratio: float = {t1m1_c1_m_c_1, ..., t1m4_c1_m_c_1,
2   ..., t5m2_c1_m_c_1, ..., t6m1_c1_m_c_1}.average()
```

By carrying out the above calculations for each column from `detailed_time_lines` we receive modules executions times ratios:

```
1   >>> times_ratios = getModulesExecutionTimesRatios(
2   test_tasks_executions, [c1, c2])
3   >>> print(times_ratios)
4   {
5   (c1, 'dl_0'): [1.41, 1.0, 1.02, 1.20],
6   ...,
7   (c2, 'dl_0'): [1.32, 1.55, 1.28, 2.04],
8   ...,
9   (c2, 'dl_15'): [1.32, 1.55, 1.28, 2.04],
10  ...,
11  (c2, 'dl_767'): [1.0, 1.08, 1.1, 1.0],
12  }
```

6. `getPricingMachinesRatios(detailes_time_lines: Dict[str: List[Tuple[Module, Machine]]], [c1, c2]) → Dict[str: List[float]]` - for each column from `detailed_time_lines` we count the ratios of machines pricing:

```
1   >>> pricing_ratios = getPricingMachinesRatios(
2   test_tasks_executions, detailes_time_lines, [c1, c2])
3   >>> print(pricing_ratios)
4   {
5   (c1, 'dl_0'): [1.0, 1.0, 1.0, 1.0],
6   ...,
7   (c2, 'dl_0'): [1.0, 1.11, 1.08, 1.00],
8   ...,
9   (c2, 'dl_15'): [1.25, 1.25, 1.25, 1.25],
10  ...,
11  (c2, 'dl_767'): [2.08, 2.08, 2.08, 2.08],
12  }
```

7. `getTimePricingRatios(detailes_time_lines: Dict[str: List[Tuple[Module, Machine]]], [c1, c2]) → Dict[str: float]` - this function return the ending result of the algorithm as a time-pricing ratios which are composed of two parts:

- (a) Time part, calculated using the formula below:

$$t = detailed\_time\_line.sum() * (1 + app-1.AIFR(task.input\_data) * cluster\_IFR) \quad (4)$$

- (b) Price part:

$$p = (detailed\_time\_line * pricing\_ratios).sum() \quad (5)$$

We normalize this part through the all details\_lines (is it necessary?) and multiply the parts by the user price-ratio parameter respectively:

$$tpr = time\_price\_ratio * t + (1 - time\_price\_ratio) * p \quad (6)$$

```

1  >>> time_pricing_ratios = getTimePricingRatios(
2      internal_flows_ratios,
3      times_ratios,
4      pricing_ratios,
5      time_price_ratio
6  )
7  >>> print(pricing_ratios)
8  {
9      (c1, 'dl_0'): 1.37,
10     ...,
11     (c2, 'dl_0'): 1.39,
12     ...,
13     (c2, 'dl_15'): 1.21,
14     ...,
15     (c2, 'dl_767'): 1.0,
16 }
17

```

It needs to be highlighted that we sort clusters for a JobBatch. JobBatch is a part of a Task and if we consider a Task with a few JobBatches we have to include data transfer time between cluster in the sorting algorithm. It will be realized through introducing the data transfer time coefficient (DTTC). For instance, if a previous JobBatch was executed at Cluster 'A', then actual JobBatch will have the DTTC equal to one for Cluster 'A' and equal to more than one to Cluster 'B'. Time-pricing ratios will be multiply by the corresponding DTTC and thus we will get the final time-pricing ratios. The smaller time-pricing ratio a time line have the better it suits the user expectations.

## 4 Algorithms tests

## 5 Conclusions