

Sort clusters algorithms

December 21, 2020

December 21, 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 Resources:
2     MEMORY: float # [GB]
3     vCPU: int # [quantity of virtual CPU's]
4
```

The sorting algorithm assume that every virtual Central Process Unit (vCPU) is equal in terms of computation speed. This is a significant generalization because actual CPU can have different clock speeds (most of them have 2-3 GHz, but it could have even more than 8 GHz[2]), have different manufacturer and different specification.

```
1 class Machine:
2     RESOURCES: resources
3     PRICE: float # [EURO/h]
4
```

Listing 1: 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
```

```

5 class Cluster:
6     ID: int
7     MACHINES: Dict[str: Machine]
8

```

Listing 2: Cluster data structure

2.2 Application-based structures

```

1 class Data:
2     TYPE: Enum['SET_OF_IMAGES', 'CSV_FILE', ...]
3
4
5 class DataSpec(Data):
6     QUANTITY: int # The number of data parts, if it is a film it
7                   # should be 1, if it is a set of photos, it should be the number
8                   # of photos
9     SIZES: [float] # size of each data part [GB]
10
11
12 class Module:
13     NAME: str
14     REQUIRED_RESOURCES: Resources
15     INPUT_DATA: Data
16

```

Listing 3: Module data structure

```

1 class JobBatch:
2     CLUSTER_ID: int
3     MODULES: Dict[str: Module]
4     INPUT_DATA: Data
5     OUTPUT_DATA: Data
6     CAL_DIAGRAM: str # CAL diagram in .yaml format
7
8
9 class Task:
10     JOB_BATCHES: List[JobBatch]
11     INPUT_DATA: DataSpec
12     RESOURCES_RESERVATION_RANGE: Tuple[Resources, Resources]
13
14
15 class TaskExecutionTime:
16     TASK_ID: int
17     TIME: float # [s]
18

```

Listing 4: TaskExecution data structure

```

1 class ModuleSpec:
2     MODULE_ID: int
3     INPUT_DATA: DataSpec
4

```

Listing 5: ModuleSpec data structure

```

1 class ModuleExecutionTime:
2     MODULE_SPEC_ID: int
3     MACHINE: Machine
4     TIME: float # [s]
5

```

Listing 6: ModuleExecutionTime data structure

3 Algorithm implemetation

3.1 Input and output data

In this subsection we will describe an input and output data for the sorting algorithm. Moreover, we will provide an example input data that will be used to make the algorithm description more clear.

As a input data for the algorithm we will use:

1. *film_processor_job_batch* (Fig. 1): JobBatch - the job batch that we want to execute,
2. *matching_clusters*: List[Cluster] - list of clusters that meet resources reservation range of the *film_processor_job_batch*
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 list of a pairs Module to Machine for each cluster. Additionally we receive the instruction (configuration) of how to run each module (the internal modules parrallelization) and when (parallelization between modules). You can find more about modules parrallelization in the algorithm scheme indtroduced further.

If there are appropriate data, the algorithm besides, it sorts clusters, predict the execution time of a job batch. Otherwise, we just receive the relative (not real) executions times. The same situation we have with overall price. We should provide specific data to receive the real price estimation. The example output can look like this (Machine, Cluster and Module simplified data structures are shown in listening 1, 2 and 3):

```

1 [
2     (
3         c_1: Cluster,
4         [ # These are the module to machine bindings with modules
5           parrallelization configuration
6           (m_1: Module, M_1: Machine, n: float, block: float),
7           (m_2: Module, M_1: Machine, n: float, block: float),
8           (m_3: Module, M_4: Machine, n: float, block: float),
9           ...
10        ],
11        overall_time: float,
12        is_time_real: bool,
13        overall_price: float,
14        is_price_real: bool,
15    ),

```

```

15 (
16   ...
17 ),
18   ...
19 ]
20

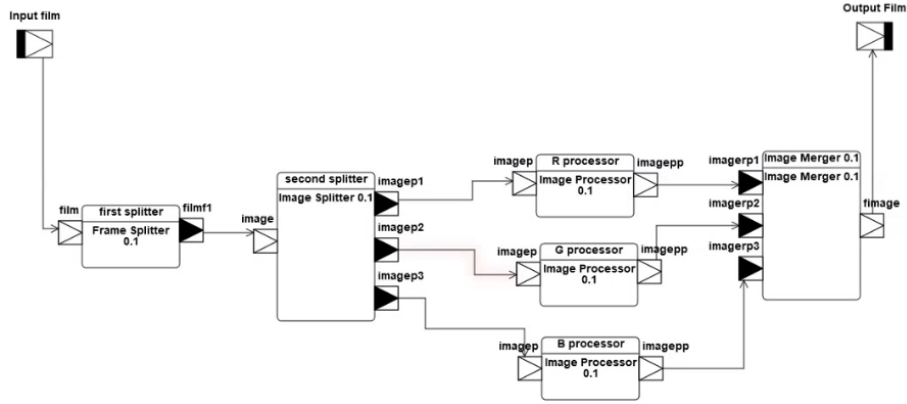
```

Listing 7: Example output of the sorting algorithm

3.2 Algorithms utils

Let's get through our algorithms with a job batch as follows:

Figure 1: *film_processor_job_batch* implementation with CAL language



```

1  film_processor_job_batch = JobBatch(
2    MODULES = {
3      'first splitter': m1,
4      'second splittter': m2,
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 )

```

An example Task structure is listen below:

```

1  film_processor_task = Task(

```

```

2 JOB_BATCHES = [film_processor_job_batch],
3 INPUT_DATA: DataSpec(film_processor_job_batch.INPUT_DATA) (
4     QUANTITY = 1
5     SIZES = [1.2]
6 ),
7 RESOURCES_RESERVATION_RANGE: (
8     Resources(
9         vCPU = 1,
10        MEMORY = 1,
11    ),
12    Resources(
13        vCPU = 32,
14        MEMORY = 32,
15    )
16 )
17 )
18

```

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

- *gcp_cluster*[1]
- *aws_cluster*[3]
- *tartu_cluster*[4]

In the listening below we can find details for each cluster (this are the original machines and prices from the given providers, limited to the most popular machines for the sake of simplicity in explaining the sorting algorithm):

```

1 gcp_cluster: Cluster = Cluster(
2     ID: 1,
3     MACHINES: {
4         'e2-standard-2': Machine = (
5             RESOURCES: {
6                 MEMORY = 8
7                 vCPU = 2
8             }
9             PRICE = 0.067
10        ),
11        'e2-standard-4': Machine = (
12            RESOURCES: {
13                MEMORY = 16
14                vCPU = 4
15            }
16            PRICE = 0.13
17        ),
18        'e2-highcpu-8': Machine = (
19            RESOURCES: {
20                MEMORY = 8
21                vCPU = 8
22            }
23            PRICE = 0.20
24        )
25    },
26 ),
27 aws_cluster: Cluster = Cluster(

```

```

28 ID: 2,
29 MACHINES: {
30     'm5.large': Machine = (
31         RESOURCES: {
32             MEMORY = 2
33             vCPU = 4
34         },
35         PRICE = 0.135
36     ),
37     'm5.xlarge': Machine = (
38         RESOURCES: {
39             MEMORY = 4
40             vCPU = 16
41         },
42         PRICE = 0.23
43     ),
44     'r5.xlarge': Machine = (
45         RESOURCES: {
46             MEMORY = 4
47             vCPU = 32
48         },
49         PRICE = 0.304
50     )
51 },
52 ),
53 tartu_cluster: Cluster = Cluster(
54     ID: 3,
55     MACHINES: {
56         'small': Machine = (
57             RESOURCES: {
58                 MEMORY = 6
59                 vCPU = 1
60             },
61             PRICE = 0.014
62         ),
63         'megium': Machine = (
64             RESOURCES: {
65                 MEMORY = 6
66                 vCPU = 4
67             },
68             PRICE = 0.064
69         ),
70         'large': Machine = (
71             RESOURCES: {
72                 MEMORY = 16
73                 vCPU = 8
74             },
75             PRICE = 0.112
76         ),
77     },
78 )
79

```

3.3 Algorithm scheme

Let's introduce the following assumptions:

1. Each instance of any vCPU is equal in terms of computation speed, regardless of the cluster it belongs to.
2. The output is a list of sorted clusters with task execution details (parallelization configuration).
3. Algorithm estimates the relative (in relation to other clusters) costs and times of task execution. It can provide real time and price estimation but the appropriate amount of specific historical data is required.

Algorithm consist of the following steps:

1. Construct the modules data specifications. Running the JobBatch, we have only the specification of input data for the first module. Subsequent modules take the data from the previous ones and we can use any method to predict (or establish from the module specification) the output data specification. In our example we assume that all of the data specifications have the same size and quantity (the size and quantity of the JobBatch input data) or the quantity and size will be reduce based on input data type of specific module.

```

1      >>> getModulesDataSpec(film_processor_job_batch)
2      [
3          'first_splitter': DataSpec(m1.INPUT_DATA) (
4              QUANTITY = 1,
5              SIZES = [1.2],
6          ),
7          'second_splitter': DataSpec(m2.INPUT_DATA) (
8              QUANTITY = 1,
9              SIZES = [1.2],
10         ),
11         'R processor': DataSpec(m3.INPUT_DATA) (
12             QUANTITY = 1,
13             SIZES = [1.2],
14         ),
15         'G processor': DataSpec(m3.INPUT_DATA) (
16             QUANTITY = 1,
17             SIZES = [1.2],
18         ),
19         'B processor': DataSpec(m3.INPUT_DATA) (
20             QUANTITY = 1,
21             SIZES = [1.2],
22         ),
23         'Image Merger 0.1': DataSpec(m4.INPUT_DATA) (
24             QUANTITY = 1,
25             SIZES = [1.2],
26         ),
27     ]
28

```

2. Return the list of module to machine pairs (with probabilities of each matching) for each cluster. It takes into the consideration resource reservation range and resources requirements for each module. The probability

can depends on how frequently the modules are assigned to each machine (the extended version of the algorithm using historical data). In this simple example we assume uniform distribution of probabilities. The function takes into the consideration only the machines from clusters for which resources are in the Task resources reservation range. In our example all of the machines at all of clusters are in the resources reservatuin range for the *film_processor_job_batch*.

```

1      {
2          1: [
3              (m1, [
4                  ('e2-standard-2', 1/3),
5                  ('e2-standard-4', 1/3),
6                  ('e2-highcpu-8', 1/3),
7              ])
8          ],
9          ...
10         (m6, [
11             ('e2-standard-2', 1/3),
12             ('e2-standard-4', 1/3),
13             ('e2-highcpu-8', 1/3),
14         ])
15     ],
16     ],
17     2: [
18         (m1, [
19             ('m5.large', 1/3),
20             ('m5.xlarge', 1/3),
21             ('r5.xlarge', 1/3),
22         ])
23     ],
24     ...
25     (m6, [
26         ('m5.large', 1/3),
27         ('m5.xlarge', 1/3),
28         ('r5.xlarge', 1/3),
29     ])
30 ],
31 ],
32 3: [
33     (m1, [
34         ('small', 1/3),
35         ('medium', 1/3),
36         ('large', 1/3),
37     ])
38 ],
39 ...
40 (m6, [
41     ('small', 1/3),
42     ('medium', 1/3),
43     ('large', 1/3),
44 ])
45 ],
46 ]
47 }
48

```


3. Estimate the module execution time for each of the pairs (module to machine) with no zero propability. The estimation of execution time for each module-to-machine pair is based on historical data of the module executions. For each module we create the simple machine learning model with the execution time as the explained variable and the following features:

- machine CPU resources where the module was executed
- data input features like the total size of the input data or the number of elements of the input data set

As we know the input data and the resources of the machine on which the module will be executed, we can ask our model about the execution time estimation. Here we have an example output for the first and the last pair of cluster 1 (we have replaced the machine name with it price to make futher price calculations possible):

```

1      {
2        1: [
3          (m1, [
4            (0.067, 1/3, 32.54),
5            (0.13, 1/3, 28.44),
6            (0.20, 1/3, 21.11),
7          ]
8        ),
9        ...
10       (m6, [
11         (0.067, 1/3, 122.11),
12         (0.13, 1/3, 111.47),
13         (0.20, 1/3, 91.26),
14       ]
15     )
16   ],
17   ...
18 }
19
```

4. ^{*1} If we have enough computational power assured by the resources reservation range and the module can be run in parallel it can decrease the execution time by n times, where n is the amount of machine instances executing the module. Figure 2 shows the example internal modules parallelization with n equal to 1, 3, 4 and 1 respectively for module 1, 2, 3 and 4.

Moreover, different modules can be executed in the same time if the application specification allows to do that. Let's assume that the *first_splitter* (which splits film into frames) can be executed part by part in the same time with *second_splitter* (which takes the image and return the 3 images, one for each RGB channel) and the *R/G/B_processor* as well. With the

¹Steps marked with * sign are the additional steps to provide the better estimations of execution time and price but are not necessary to perform clusters sorting.

assumption described above the *film_processing_task*'s execution time line can look like it is shown on the figure 3.

Since we estimated the execution times for each module in the previous step, now we can find all possible parallelization configuration for the job batch and further search for the best one.

In our example we simplify this step by providing only the single one configuration where there is no any parallelization at all:

```

1  {
2    1: [
3      "c1_conf1": [
4        (m1, [
5          (0.067, 1/3, 32.54),
6          (0.13, 1/3, 28.44),
7          (0.20, 1/3, 21.11),
8        ], 1, 0
9      ),
10     ...
11     (m6, [
12       (0.067, 1/3, 122.11),
13       (0.13, 1/3, 111.47),
14       (0.20, 1/3, 91.26),
15     ], 1, 5
16   )
17 ],
18 ],
19 ...
20 }
21

```

The last part of each module element for a given configuration tell us which modules can be executed in the same time. One part before last is the n factor (which in our example is equal to 1 everywhere) and it will be used further to divide the modules execution time estimations and also provide information about the way of a task execution.

5. * As we can see on the Figure 2 there could be some time windows between modules executions. This windows represents the time to start the instance of machine and other preparations before module execution like start up the image or download the data. This step is sophisticated and it should not provide a large changes to final result but if we would like to have better estimations it should be carry out. It could be done by:

- Intodruce the function that returns internal flow ratio of the Job-Batch execution based on input data size and the parallelization configuration. In our example AIRF (*application interfnal flow ratio*) function can look like it is shown on the figure 4. To be more specific, it returns overall time of task execution minus sum of the maximum executions times of each modules from the task execution

time line. The function use the following formula:

$$coeff = \frac{internal_flow_time}{task_execution_time} \quad (1)$$

In the example timeline shown on the figure 2 (where modules internal parallization respectively equals 1, 3, 4, 1) the coefficient will eqaul:

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

So as a result we get the internal flow coefficient (*IFC*) in range [0, 1] the 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 3. In such a case the *IFC* will be significantly lower. Even though, it is possible to count the *IFC* for such a cases taking into the account the task execution configuration. Having the *IFC* for a task we just multiply each estimated module execution time by factor $1 + IFC$.

- Other approximation of the internal data flow time by looking at the historical data of modules executions.

6. An user responsible for the task execution is providing *time-price* ratio which is the factor prioritizing the cluster sorting. It is a float number from 0 to 1 telling the algorithm if user is willing to save more time or money.

In this step we count the best parallelization configuration for each cluster based on the *time-price* ratio. For each configuration we aggregate a time and a price for each module receiving the overall time and price of a task execution. Next we normalize the prices and times for each configuration within the cluster and use the *time-price* ratio to find the best one.

The listing below show the aggregation result:

```

1  {
2    1: [
3      (m1, (3.33, 27.36), 1, 0),
4      ...
5      (m6, (13,49, 108,28), 1, 5),
6    ],
7    ...
8  }
9
```

In this step we had ready the estimates of execution price and time for each module on each cluster.

7. * 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. *DTTC* can be based on historical data of jobs executions that run different batches on different clusters. As this step is enough sophisticated to perform with more complex scenarios, it is omitted by now.

8. The final step is counting the overall time and price for each cluster and then sort them by using the time-price ratio. We got the following overall times and prices with the task execution configurations:

```

1  {
2      1: (41.22, 287.17, "c1_conf1"),
3      2: (61.81, 221.18, "c2_conf1"),
4      3: (51.11, 181.19, "c3_conf1"),
5  }
6

```

Having the *time-price* ratio equal to 0.6 (which means that the user is willing to save time a little more than money) we got the following final penalty for each cluster:

```

1  {
2      1: 41.22 * (1-0.6) + 287.17 * 0.6 = 188,79,
3      2: 61.81 * 0.4 + 221.18 * 0.6 = 157,42,
4      3: 51.11 * 0.4 + 181.19 * 0.6 = 129,16,
5  }
6

```

The smaller penalty is, the more suitable the cluster is for given task execution and *time-price* ratio. The final list of sorted clusters:

```

1  [
2      (3, 129,16, "c3_conf1"),
3      (2, 157,42, "c2_conf1"),
4      (1, 188,79, "c1_conf1"),
5  ]
6

```

If the probabilities of matching module to machine have uniform distribution (this happen when we do not have historical data of modules executions) and all of the machines resources are in the task resources reservation range, and we do not use any parallelization, the algorithm is just counting the average price of a machine for each cluster and that what the sorting is based on.

Figure 2: Example of the *file_processor_task* timeline with internal module parallelization.

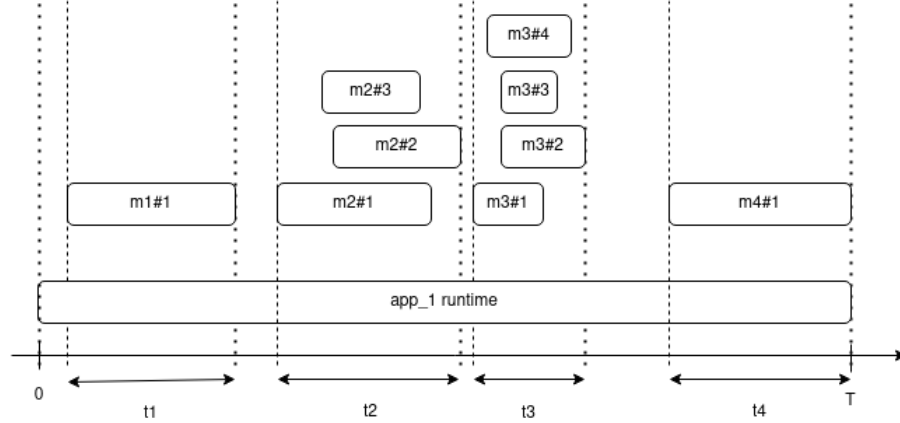
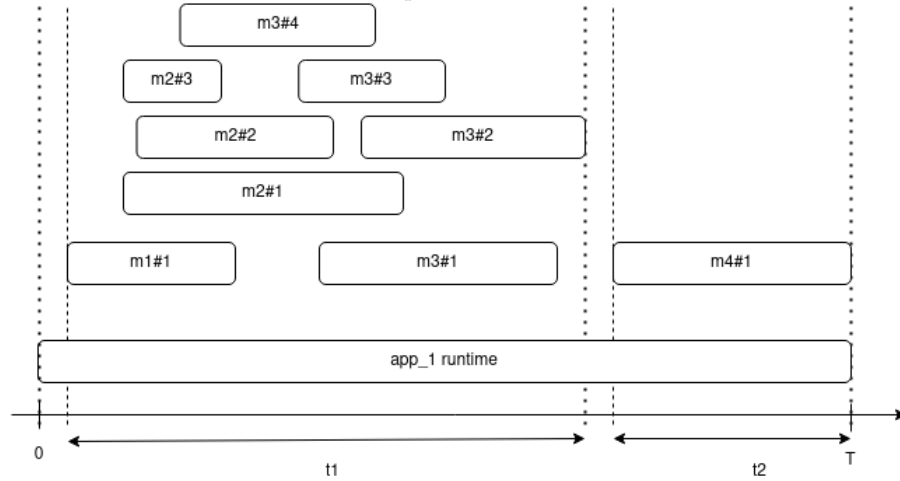


Figure 3: Example of the *file_processor_task* timeline with internal module parallelizations and between modules parallelization as well.



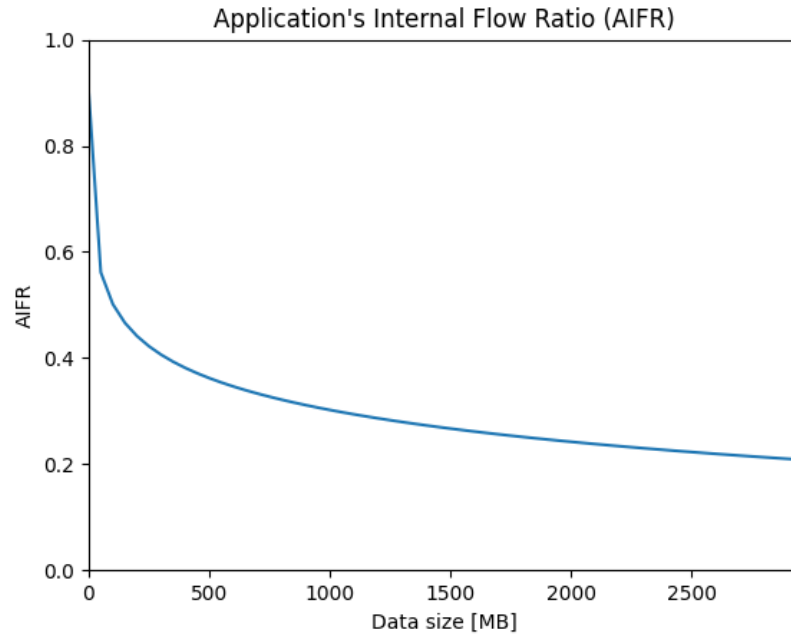
4 Algorithms tests

5 Conclusions

References

- [1] Google: Google Cloud Platform - pricing, 03.10.2020, <https://cloud.google.com/compute/all-pricing>

Figure 4: *file_processor_task* AIFR



- [2] asus.com: 8.79GHz FX-8350 Is The Fastest Ever CPU, 03.10.2020,
<https://rog.asus.com/articles/crosshair-motherboards/8-79ghz-fx-8350-is-the-fastest-ever>
- [3] Amazon: Amazon Web Services - pricing, 04.10.2020,
<https://aws.amazon.com/emr/pricing/>
- [4] University of Tartu, HPC - pricing, 04.10.2020,
<https://hpc.ut.ee/en/prices/>