# SLOT SELECTION AND CO-ALLOCATION ALGORITHMS FOR ECONOMIC SCHEDULING IN DISTRIBUTED COMPUTING

**Victor Toporkov, Dmitry Yemelyanov**

National Research University "MPEI"
Moscow / Russia
*E-mail: ToporkovVV@mpei.ru, yemelyanov.dmitry@gmail.com*

**Anna Toporkova**

National Research University Higher School of Economics
Moscow / Russia
*E-mail: atoporkova@hse.ru*

**Alexey Tselishchev**

CERN (European Organization for Nuclear Research)
Genève / Switzerland
*E-mail: Alexey.Tselishchev@cern.ch*

## Abstract

In this work, we introduce slot selection and co-allocation algorithms for parallel jobs in distributed computing with non-dedicated resources. A single slot is a time span that can be assigned to a task, which is a part of a job. The job launch requires a co-allocation of a specified number of slots starting synchronously. The challenge is that slots associated with different CPU nodes of distributed computational environments may have arbitrary start and finish points that do not match. Some existing algorithms assign a job to the first set of slots matching the resource request without any optimization (the first fit type), while other algorithms are based on an exhaustive search. In this paper, algorithms for effective slot selection of linear complexity on an available slots number are studied and compared with known approaches. The novelty of the proposed approach consists of allocating alternative sets of slots. It provides possibilities to optimize job scheduling.

*Keywords -* Distributed computing; economic scheduling; resource management; slot; job; batch.

## 1    INTRODUCTION

Economic mechanisms are used to solve tasks like resource management and scheduling of jobs in a transparent and efficient way in distributed environments such as utility Grid and cloud computing [1, 2]. A resource broker model is decentralized, well-scalable and application-specific [1-4]. It has two parties: node owners and brokers representing users. The simultaneous satisfaction of various application optimization criteria submitted by independent users is unreachable in essence and also can deteriorate such quality of service rates as total execution time of a sequence of jobs or overall resource utilization. Another model is related to virtual organizations (VO) [5-7] with central schedulers providing job-flow level scheduling and optimization. VOs naturally restrict the scalability, but uniform rules for allocation and consumption of resources make it possible to improve the efficiency of resource usage and find a trade-off between contradictory interests of different participants.

In [6], we have proposed a hierarchical model of resource management system which is functioning within a VO. Resource management is implemented with a structure consisting of a metascheduler and subordinate job schedulers that interact with batch job processing systems. The significant difference between the approach proposed in [6] and well-known scheduling solutions for distributed environments such as Grids [1-5, 8, 9], e.g., gLite Workload Management System [8], where Condor [10] is used as a scheduling module, is the fact that the scheduling strategy is formed on a basis of efficiency criteria. They allows to reflect economic principles of resource allocation by using relevant

cost functions and solving a load balance problem for heterogeneous processor nodes. At the same time the inner structure of the job is taken into account when the resulting schedule is formed. The metascheduler [5-7, 10] implements the economic policy of a VO based on local CPU schedules. The schedules are defined as sets of slots coming from local resource managers or schedulers in the node domains. During each scheduling cycle the sets of available slots are updated based on the information from local resource managers. Thus, during every cycle of the job batch scheduling [6] two problems have to be solved: 1) selecting alternative set of slots (alternatives) that meet the requirements (resource, time, and cost); 2) choosing a slot combination that would be the efficient or optimal in terms of the whole job batch execution in the current cycle of scheduling. To implement this scheduling scheme, first of all, one needs to propose the algorithm for finding sets of alternative executions. An optimization technique for the second phase of this scheduling scheme was proposed in [6, 7]. First fit slot selection algorithms assign any job to the first set of slots matching the resource request conditions, while other algorithms use an exhaustive search. Moab scheduler [11, 12] implements backfilling algorithm and during a slot window search does not take into account any additive constraints such as the minimum required storage volume or the maximum allowed total allocation cost. Moreover, backfilling does not support environments with non-dedicated resources and its execution time grows substantially with increase of the slot number. Assuming that every CPU node has at least one local job scheduled, the backfilling algorithm has quadratic complexity in the slot number. In [4] heuristic algorithms for slot selection, based on user-defined utility functions, are introduced. NWIRE system [4] performs a slot window allocation based on the user defined efficiency criterion under the maximum total execution cost constraint. However, the optimization occurs only on the stage of the best found offer selection. In our previous works [13-15], two algorithms for slot selection AMP and ALP that feature linear complexity $O(m)$, where $m$ is the number of available time-slots, were proposed. Both algorithms perform the search of a first fitting window without any optimization. AMP (Algorithm based on Maximal job Price), performing slot selection based on the maximum slot window cost, proved the advantage over ALP (Algorithm based on Local Price of slots) when applied to the above mentioned scheduling scheme. However, in order to accommodate an end user's job execution requirements, there is a need for a more precise slot selection algorithm of linear complexity to exploit during the first stage of the proposed scheduling scheme and to consider various user demands along with the VO resource management policy.

In this paper, we propose algorithms for effective slot selection based on user defined criteria that feature linear complexity on the number of the available slots during the job batch scheduling cycle. The novelty of the proposed approached consists in allocating a number of alternative sets of slots (alternatives). The proposed algorithms can be used for both homogeneous and heterogeneous resources. The paper is organized as follows. Section 2 introduces a general scheme for searching alternative slot sets that are effective by the specified criteria. Then four implementations are proposed and considered. Section 3 contains simulation results for comparison of proposed and known algorithms. Section 4 summarizes the paper and describes further research topics.

## 2  GENERAL SCHEME AND SLOT SELECTION ALGORITHMS

In this section we introduce a problem statement, consider a general scheme of an **A**lgorithm searching for **E**xtreme **P**erformance (AEP) and its implementation examples.

### 2.1  Problem statement and AEP scheme

The launch of any job requires co-allocation of a specified number of slots, as well as in the backfilling conservative variation [11, 12]. A single *slot* is a time span that can be assigned to *a task*, which is a part of *a job*. According to the resource request, it is required to find a window $W$ with the following description: $n$ concurrent time-slots providing the resource performance rate and the maximal resource price per time unit $F$ should be reserved for a time span $T_i$. The task is to scan a list of $m$

available slots and to select a window $W$ of $n$ parallel slots with a length of the required resource reservation time. The total window cost is calculated as a sum of an individual usage cost of the selected slots. In addition, one can define a criterion on which the best matching window alternative is chosen. This can be a criterion $crW$ for a minimum cost, a minimum execution runtime or, for example, a minimum energy consumption.

Consider as an example the problem of selecting a window of size $n$ with a total cost no more than $S$ from the list of $m > n$ slots (in the case, when $m = n$ the selection is trivial). The maximal job budget is counted as $S = FT_i n$. The current extended window consists of $m$ slots $s_1, s_2, ..., s_m$. The cost of using each of the slots according to their required time length is: $c_1, c_2, ..., c_m$. Each slot has a numeric characteristic $z_i$ in accordance to $crW$. The total value of these characteristics should be minimized in the resulting window. Then the problem could be formulated as follows:

$$a_1 z_1 + a_2 z_2 + ... + a_m z_m \to \min, \qquad a_1 c_1 + a_2 c_2 + ... + a_m c_m \le S, \qquad a_1 + a_2 + ... + a_m = n,$$

$a_r \in \{0, 1\}$, $r = 1, ..., m$. Additional restrictions can be added, for example, considering the specified value of deadline. Finding the coefficients $a_1, a_2, ..., a_m$ each of which takes integer values 0 or 1 (and the total number of "1" values is equal to $n$), determine the window with the specified criteria $crW$ extreme value. The algorithm parses a ranged list of all available slots subsequently for all the batch jobs. Existing slot selection algorithms assign a job to the first set of slots matching the resource request conditions or use an exhaustive search. AEP is free of the obvious disadvantages of the exhaustive search and has linear complexity on the number of the slots available in the current scheduling cycle. AEP can be compared to the algorithm of min/max value search in an array of flat values. The effective on the specified criterion window of size $n$ is selected from this $m$ slots and compared with the results in the previous steps. By the end of the slot list the only solution with the best criterion $crW$ value will be selected. The algorithm proposed is processing a list of all slots available during the scheduling interval ordered by a non-decreasing start time. This condition is required for a single sequential slot list scan and algorithm linear complexity on the number $m$ of slots.

The scheme for an effective window search by the specified criteria can be represented as follows:

```
/* Job – Batch job for which the search is performed ;
** windowSlots – a set (list) of slots representing the window;*/
slotList = orderSystemSlotsByStartTime();
for(i=0; i< slotList.size; i++){
        nextSlot = slotList[i];
 if(!properHardwareAndSoftware(nextSlot))
     continue;    // The slot does not meet the requirements
 windowSlotList.add(nextSlot);
 windowStart = nextSlot.startTime;
 for(j=0; j<windowSlots.size; j++){
  wSlot = windowSlots[j];
  minLength = wSlot.Resource.getTime(Job);
  if((wSlot.EndTime – windowStart) < minLength)
     windowSlots.remove(wSlot);
 }
 if(windowSlots.size >= Job.m){
   curWindow = getBestWindow(windowSlots);
   crW = getCriterion(curWindow);
   if(crW > maxCriterion){
      maxCriterion = crW;
     bestWindow = curWindow;
```

```
        }
    }
}
```

Finally, a variable `bestWindow` will contain an effective window by the given criterion.

## 2.2   AEP implementation examples

The need to choose alternative sets of slots for every batch job increases the complexity of the whole scheduling scheme. With a large number of the available slots the algorithm execution time may become inadequate. Though it is possible to mention some typical optimization tasks, based on the AEP scheme that can be solved with a relatively decreased complexity. These include problems of total job cost minimizing, total runtime minimizing, the window formation with the minimal start/finish time.

Consider the procedure for *minimizing a window start time*. The difference with the general AEP scheme is that the first suitable window will have the earliest possible start time. Indeed, if at some step *i* of the algorithm (after the *i*-th slot is added) the suitable window can be formed, then the windows, formed at the further steps will be guaranteed to have the start time that is not earlier (according to the ordered list of available slots, only slots with non-decreasing start time will be taken into consideration). Thus this procedure can be reduced to finding a set of the first $n$ parallel slots the total cost of which does not exceed the budget limit $S$. This description coincides the AMP scheme considered in previous works [13-15]. Thus AEP is naturally an extension of AMP, and AMP is the particular case of the whole AEP scheme performing only the start time optimization. Further we will use AMP abbreviation as a reference to the window start time minimization procedure.

It is easy to provide the implementation of the algorithm of finding a window with *the minimum total execution cost*. For this purpose in the AEP search scheme $n$ slots with the minimum sum cost should be chosen. If at each step of the algorithm a window with the minimum sum cost is selected, at the end the window with the best value of the criterion $crW$ will be guaranteed to have overall minimum total allocation cost at the given scheduling interval.

The task to find a window with *the minimum runtime* is more complicated. Given the nature of determining a window runtime, which is equal to the length of the longest slot (allocated on the node with the least performance level), the following algorithm may be proposed:

```
orderSlotsByCost(windowSlotList);

resultWindow = getSubList(0,n, windowSlotList);

extendWindow = getSubList(n+1,m, windowSlotList);

while(extendWindow.size > 0){

  longSlot = getLongestSlot(resultWindow);

  shortSlot = getCheapestSlot(extendWindow);

  extendWindow.remove(shortSlot);

  if((shortSlot.size < longSlot.size)&&

   (resultWindow.cost + shortSlot.cost < S)){

     resultWindow.remove(longSlot);

     resultWindow.add(shortSlot);

  }

}
```

As the result, the suitable window of the minimum time length will be formed in a variable `resultWindow.` The algorithm described consists of the consecutive attempts to substitute the longest slot in the forming window (the `resultWindow` variable) to another shorter one that will not

be too expensive. In case when it is impossible to substitute the slots without violating the constraint on the maximum window allocation cost, the current `resultWindow` configuration is declared to have the minimum runtime. Implementing this algorithm of window selection at each step of the AEP scheme allows to find a suitable window with the minimum possible runtime at the given scheduling interval.

An algorithm for finding a window with *the earliest finish time* has a similar structure and can be described using the runtime minimizing procedure presented above. Indeed, the expanded window has a start time `tStart` equal to the start time of the last added suitable slot. The minimum finish time for a window on this set of slots is (`tStart + minRuntime`), where `minRuntime` is the minimum window length. The value of `minRuntime` can be calculated similar to the runtime minimizing procedure described above. Thus, by selecting at each step of the algorithm a window with the earliest completion time, the required window will be allocated in the end of the slot list. It is worth mentioning that the all proposed AEP implementations have a linear complexity $O(m)$: algorithms "move" through the list of the $m$ available slots in the direction of non-decreasing start time without turning back or reviewing previous steps.

## 3    EXPERIMENTAL STUDIES OF SLOT SELECTION ALGORITHMS

The goal of the experiment is to examine AEP implementations: to analyze alternatives search results with different efficiency criteria, to compare the results with AMP and to estimate the possibility of using in real systems considering the algorithms' execution time.

### 3.1    Algorithms and simulation environment

For the proposed AEP efficiency analysis the following implementations were added to the simulation model [6]:

1. *AMP* – the algorithm for searching alternatives with the earliest start time. This scheme was introduced in works [13-15] and briefly described in section 2.

2. *minFinish* – the algorithm for searching alternatives with the earliest finish time.

3. *minCost* – the algorithm for searching a single alternative with the minimum total allocation cost on the scheduling interval.

4. *minRuntime* – this algorithm performs a search of a single alternative with the minimum runtime.

5. *minProcTime* – this algorithm performs a search for a single alternative with the minimum total node execution time. It is worth mentioning that this realization is simplified and does not guarantee an optimal result and only partially matches the AEP scheme, because a random window is selected.

6. *Common Stats, AMP* (further referred to as *CSA*) – the scheme for searching multiple alternatives using *AMP*. Similar to the general searching scheme [13-15], a set of suitable alternatives, disjointed by the slots, is allocated for each job. To compare the search results with the algorithms 1-5, presented above, only alternatives with the extreme value of the given criterion will be selected, so the optimization will take place at the selection process. The criteria include the minimum start time, the minimum finish time, the minimum total execution cost, the minimum runtime and the minimum processor time used.

Since the purpose of the considered algorithms is to allocate suitable alternatives, it makes sense to make the simulation apart from the whole general scheduling scheme, described in [6]. In this case, the search will be performed for a single predefined job. Thus during every single experiment a generation of a new distributed computing environment will take place while the algorithms described will perform the alternatives search for a single base job with the resource request defined in advance. A simulation framework [6, 7] was configured in order to study and analyze the presented algorithms. In each experiment a generation of the distributed environment that consists of 100 CPU nodes was performed. The performance rate for each node was generated as a random integer variable in the interval [2; 10] with a uniform distribution. The resource usage cost was formed proportionally to their performance with an element of normally distributed deviation in order to simulate a free market pricing model [1-4]. The level of the resource initial load with the local and high priority tasks at the scheduling interval [0; 600] was generated by the hyper-geometric distribution in the range from 10%

to 50% for each CPU node. Based on the generated environment the algorithms performed the search for a single initial job that required an allocation of 5 parallel slots for a 150 units of time. The maximum total execution cost according to user requirements was set to 1500. This value generally will not allow to use the most expensive (and usually the most efficient) CPU nodes. The relatively high number of the generated nodes has been chosen to allow *CSA* to find more slot alternatives. Therefore more effective alternatives could be selected for the searching results comparison based on the given criteria. The characteristics of the alternatives found (in case of *CSA* – selected) are stored for each considered algorithm during every experiment. These characteristics include alternatives' start time, finish time, runtime, total cost and total used processor time.

## 3.2  Experimental results

The results of the 5000 simulated scheduling cycles are presented in Fig. 1-4.

Consider the average *start time* for the alternatives found (and selected) by the aforementioned algorithms (Fig. 1, a).



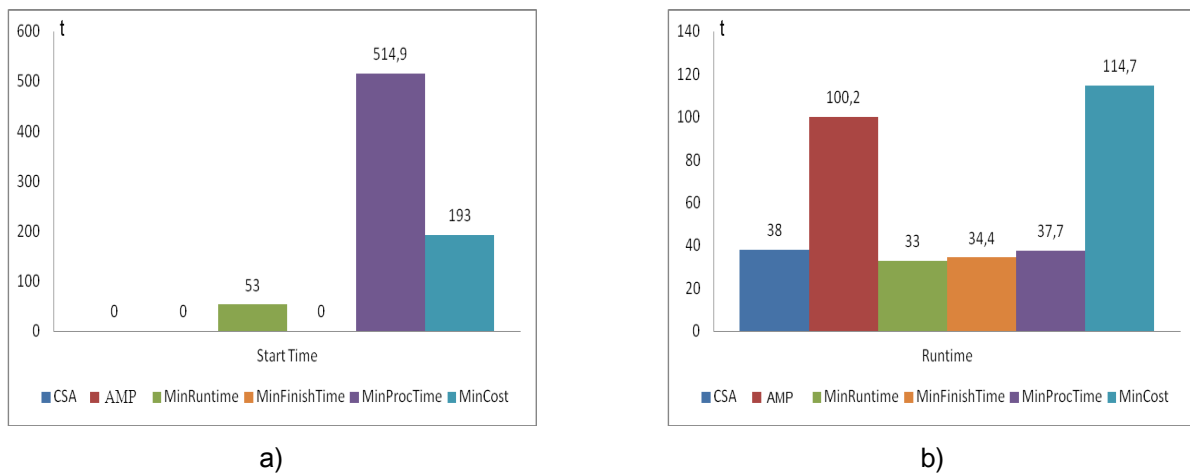a)                                                                                        b)

Fig. 1. Start time (a) and runtime (b)

*AMP*, *minFinish* and *CSA* were able to provide the earliest job start time at the beginning of the scheduling interval ($t = 0$). The result was expected for *AMP* and *CSA* (which is essentially based on the multiple runs of the *AMP* procedure) since 100 available resource nodes provide a high probability that there will be at least 5 parallel slots starting at the beginning of the interval and can form a suitable window. The fact that the *minFinish* algorithm was able to provide the same start time can be explained by the local tasks minimum length value, that is equal to 10. Indeed, the window start time at the moment $t = 10$ cannot provide the earliest finish time even with use of the most productive resources (for example the same resources allocated for the window with the minimal runtime). Average starting times of the alternatives found by *minRuntime*, *minProcTime* and *minCost* are 53, 514.9 and 193 respectively.

The *average runtime* of the alternatives found (selected) is presented in Fig. 1, b. The minimum execution runtime 33 was obviously provided with the *minRuntime* algorithm. Though, schemes *minFinish*, *minProcTime* and *CSA* provide quite comparable values: 34.4, 37.7 and 38 time units respectively that only 4.8%, 12.5% and 13.2% longer. High result for the *minFinish* algorithm can be explained by the "need" to complete the job as soon as possible with the minimum (and usually initial) start time. *MinFinish* and *minRuntime* are based on the same criterion selection procedure described in the section 2. However due to non-guaranteed availability of the most productive resources at the beginning of the scheduling interval, *minRuntime* has the advantage. Relatively long runtime was provided by *AMP* and *minCost* algorithms. For *AMP* this is explained by the selection of the first fitting (and not always effective by the given criterion) alternative, while *minCost* tries to use relatively cheap and (usually) less productive CPU nodes.

The minimum *average finish time* (Fig. 2, a) was provided by the *minFinish* algorithm – 34.4. *CSA* has the closest resulting finish time of 52.6 that is 34.6% later. The relative closeness of these values comes from the fact that other related algorithms did not intend to minimize a finish time value and were selecting windows without taking it into account. At the same time *CSA* is picking the most effective alternative among 57 (on the average) allocated at the scheduling cycle: the optimization was carried out at the selection phase. The late average finish time 307.7 is provided by the *minCost* scheme. This value can be explained not only with a relatively late average start time (see Fig. 1, a), but with a longer (compared to other approaches) execution runtime (see Fig. 1, b) due to the use of less productive resource nodes. The finish time obtained by the simplified *minProcTime* algorithm was relatively high due to the fact that a random window was selected (without any optimization) at the each step of the algorithm, though the search was performed on the whole list of available slots. With such a random selection the most effective window by the processor time criterion was near the end of the scheduling interval.



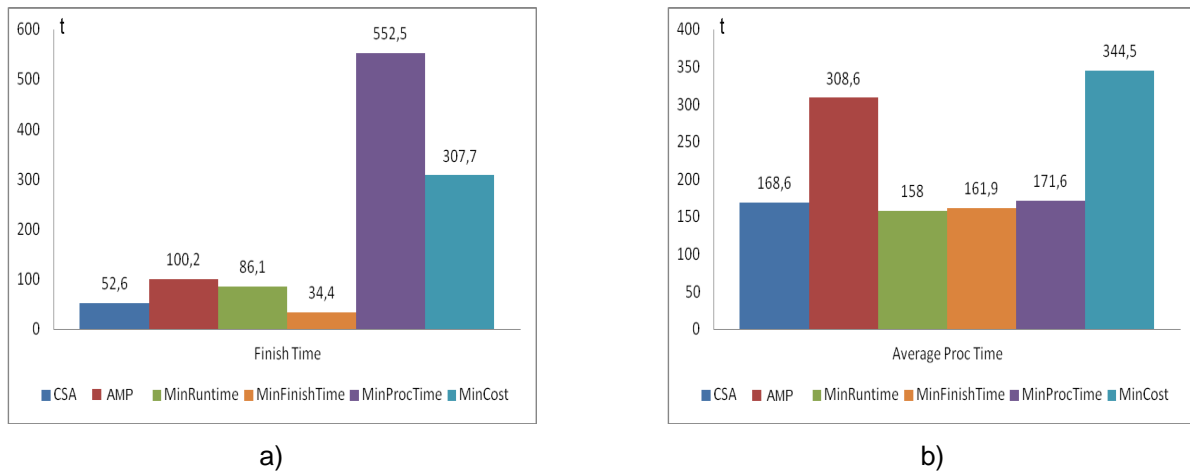a)                                                      b)

Fig. 2. Finish time (a) and CPU usage time (b)

The *average used processor time* (the sum time length of the used slots) for the algorithms considered is presented by Fig. 2, b. The minimum value was provided by *minRuntime*: 158 time units. *MinFinish*, *CSA* and *minProctime* were able to provide comparable results: 161.9, 168.6 and 171.6 respectively. It worth mentioning that although the simplified *minProcTime* scheme does not provide the best value, its only 2% less effective compared to a common *CSA* scheme, while its working time is orders of magnitude less (Tables1, 2). The most processor time consuming alternatives were obtained by *AMP* and *minCost* algorithms. Similarly to the execution runtime value, this can be explained by using a random first fitting window (in case of *AMP*) or by using less expensive hence less productive resource nodes (in case of the *minCost* algorithm), as nodes with a low performance level require more time to execute the job.

Table 1. Actual algorithms' execution time measured depending on the CPU nodes number

| CPU nodes number: | 50 | 100 | 200 | 300 | 400 |
|---|---|---|---|---|---|
| CSA:Alternatives Num | 25.9 | 57 | 128.4 | 187.3 | 252 |
| CSA per Alt | 0.33 | 0.99 | 3.16 | 6.79 | 11.83 |
| CSA | 8.5 | 56.5 | 405.2 | 1271 | 2980.9 |
| AMP | 0.3 | 0.5 | 1.1 | 1.6 | 2.2 |
| MinRuntime | 3.2 | 12 | 45.5 | 97.2 | 169.2 |
| MinFinishTime | 3.2 | 12 | 45.1 | 96.9 | 169 |

| | | | | | |
|---|---|---|---|---|---|
| MinProcTime | 1.5 | 5.2 | 19.4 | 42.1 | 74.1 |
| MinCost | 1.7 | 6.3 | 23.6 | 52.3 | 91.5 |

Finally, consider the *total job execution cost* (Fig. 3). The *minCost* algorithm has a big advantage over the other algorithms presented: it was able to provide the total cost of 1027.3 (note that the total cost limit was set by the customer at 1500). Alternatives found with other considered algorithms have approximately the same execution cost. Thus, the cheapest alternatives found by *CSA* have the average total execution cost equal to 1352, that is 24% more expensive compared to the result of the *minCost* scheme, while alternatives found by *minRuntime* (the most expensive ones) are 29.9% more expensive.

The important factor is a complexity and an actual working time of the algorithms under consideration especially with the assumption of the algorithms' repeated use during the first stage of the scheduling interval. In the description of the AEP general scheme it was mentioned that the algorithm has a linear complexity on the number of the available slots and a quadratic complexity with a respect to the number of CPU nodes. Table 1 shows the actual algorithms' execution time in milliseconds measured depending on the number of CPU nodes. The simulation was performed on a regular PC workstation with Intel Core i3 (2 cores by 2.93 GHz), 3GB RAM on JRE 1.6, and 1000 separate experiments were simulated for each value of the processor nodes numbers {50, 100, 200, 300, 400}.
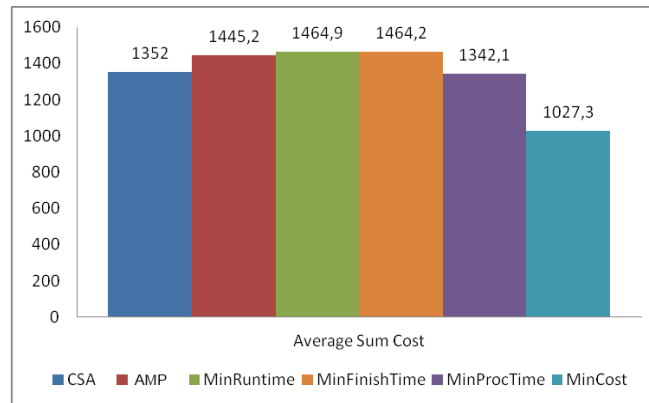


Fig. 3. Job execution cost

The simulation parameters and assumptions were the same as described in section 3.1, apart from the number of used CPU nodes. A row "*CSA: Alternatives Num*" represents an average number of alternatives found by *CSA* during the single experiment simulation (note that *CSA* is based on multiple runs of AMP algorithm). A row "*CSA per Alt*" represents an average the *CSA* algorithm working time in recalculation for one alternative.

The *CSA* scheme has the longest working time that on the average almost reaches 3 seconds when 400 nodes are available. Besides this time has a near cubic increasing trend with a respect to the nodes number.

This trend can be explained by the addition of two following factors:

1) a linear increase of the alternatives number found by *CSA* at each experiment (which makes sense: the linear increase of the available nodes number leads to the proportional increase in the available node processor time; this makes it possible to find (proportionally) more alternatives);

2) a near quadratic complexity of the *AMP* algorithm with a respect to the nodes number, which is used to find single alternatives in *CSA*. Even more complication is added by the need of "cutting" a suitable windows from the list of the available slots.

Other considered algorithms will be able to perform a much faster search. The average working time of *minRuntime*, *minProcTime* and *minCost* proves their (at most) quadratic complexity on the number of CPU nodes. The *AMP*'s execution time shows even near linear complexity because with a relatively large number of free available resources it was usually able to find a window at the beginning of the scheduling interval (see Fig. 1, a) without the full slot list scan. Fig. 4, a clearly presents the average

working duration of considering algorithms depending on the number of available CPU nodes (the values were taken from Table 1). (The *CSA* curve is not represented as its working time is incomparably longer than AEP like algorithms.)

Table 2 contains the algorithms' working time in milliseconds measured depending on the scheduling interval length.

Overall 1000 single experiments were conducted for each value of the interval length {600, 1200, 1800, 2400, 3000, 3600} and for each considered algorithm an average working time was obtained. The experiment simulation parameters and assumptions were the same as described earlier in this section, apart from the scheduling interval length. A number of CPU nodes was set to 100. Similarly to the previous experiment,

*CSA* had the longest working time (about 2.5 seconds with scheduling interval length equal to 3600 model time units), which is mainly caused by the relatively large number of the formed execution alternatives (on the average more than 400 alternatives on a 3600 interval length). When analyzing the presented values it is easy to ensure that all proposed algorithms have a linear complexity with the respect to the length of the scheduling interval and, hence, to the number of the available slots (Fig. 4, b).



a)                                                                                          b)
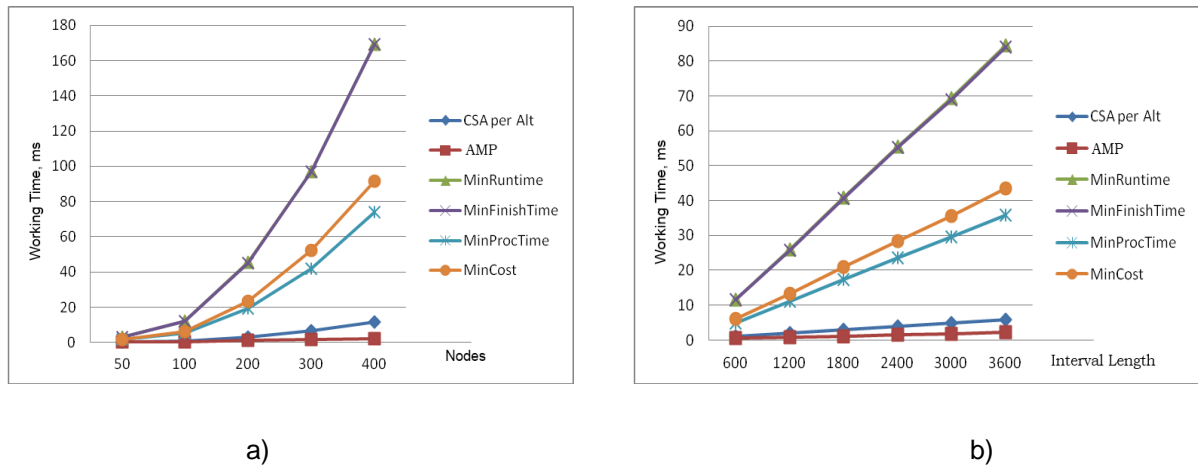
Fig. 4. Working time depending on the available CPU nodes number (a) and on the scheduling interval length (b)

Table 2. Algorithms' working time measured depending on the scheduling interval length

| Scheduling interval length: | 600 | 1200 | 1800 | 2400 | 3000 | 3600 |
|---|---|---|---|---|---|---|
| Number of slots | 472.6 | 779.4 | 1092 | 1405.1 | 1718.8 | 2030.6 |
| CSA: Alternatives Num | 57 | 125.4 | 196.2 | 269.8 | 339.7 | 412.5 |
| CSA per Alt | 0.95 | 1.91 | 2.88 | 3.88 | 4.87 | 5.88 |
| CSA | 54.2 | 239.8 | 565.7 | 1045.7 | 1650.5 | 2424.4 |
| AMP | 0.5 | 0.82 | 1.1 | 1.44 | 1.79 | 2.14 |
| MinRuntime | 11.7 | 26 | 40.9 | 55.5 | 69.4 | 84.6 |
| MinFinishTime | 11.6 | 25.7 | 40.6 | 55.3 | 69 | 84.1 |
| MinProcTime | 5 | 11.1 | 17.4 | 23.5 | 29.5 | 35.8 |
| MinCost | 6.1 | 13.4 | 20.9 | 28.5 | 35.7 | 43.5 |

## 4    CONCLUSIONS AND FUTURE WORK

In this work, we address the problem of slot selection and co-allocation for parallel jobs in distributed computing with non-dedicated resources. For this purpose AMP and AEP approaches were proposed and considered. Specific AEP scheme implementations with a reduced over a general scheme complexity were proposed and considered. Each of the algorithms possesses a linear complexity on a total available slots number and a quadratic complexity on a CPU nodes number. The advantage of AEP-based algorithms over the general *CSA* scheme was shown for each of considered criteria: start time, finish time, runtime, CPU usage time and total cost. In our further work we will refine resource co-allocation algorithms in order to integrate them with scalable co-scheduling strategies [6, 7]. Future research will be focused on  further AEP based algorithms research and its integration with the whole batch scheduling approach, and mainly on its influence on job-flows execution efficiency.

## ACKNOWLEDGEMENTS

## References

[1]  S.K. Garg, R. Buyya, H.J. Siegel, Scheduling Parallel Applications on Utility Grids: Time and Cost Trade-off Management. Proc of ACSC 2009, Wellington, New Zealand (2009) 151-159.

[2]  S.K. Garg, C.S. Yeo, A. Anandasivam, R.Buyya, Environment-conscious Scheduling of HPC Applications on Distributed Cloud-oriented Data Centers. J. of Parallel and Distributed Computing. 71 (6) (2011). 732-749.

[3]  R. Buyya, D. Abramson, J. Giddy, Economic Models for Resource Management and Scheduling in Grid computing. J. of Concurrency and Computation: Practice and Experience. 14(5) (2002) 1507–1542.

[4]  C. Ernemann, V. Hamscher, R. Yahyapour, Economic Scheduling in Grid Computing. Proc. of the 8th Job Scheduling Strategies for Parallel Processing. Eds D.G. Feitelson, L. Rudolph, U. Schwiegelshohn. Heidelberg: Springer, LNCS. 2537 (2002) 128-152.

[5]  K. Kurowski, J. Nabrzyski, A. Oleksiak et al., Multicriteria Aspects of Grid Resource Management. Grid resource management. State of the art and future trends. Eds J. Nabrzyski, J.M. Schopf and J. Weglarz. Kluwer Acad. Publ. (2003) 271–293.

[6]  V. Toporkov, A. Tselishchev, D. Yemelyanov, A. Bobchenkov, Composite Scheduling Strategies in Distributed Computing with Non-dedicated Resources. Procedia Computer Science. Elsevier. 9 (2012) 176-185.

[7]  V. Toporkov, A. Tselishchev, D. Yemelyanov, A. Bobchenkov, Dependable Strategies for Job-flows Dispatching and Scheduling in Virtual Organizations of Distributed Computing Environments. Complex Systems and Dependability. Berlin, Heidelberg: Springer-Verlag, AISC. 170 (2012) 240-255.

[8]  M. Cecchi, F. Capannini, A. Dorigo et al., The gLite Workload Management System. J. Phys.: Conf. Ser. 219 (6) (2010) 062039.

[9]  J. Yu, R. Buyya, K. Ramamohanarao, Workflow Scheduling Algorithms for Grid Computing. Metaheuristics for Scheduling in Distributed Computing  Environments, Studies in Computational Intelligence. 146. Springer-Verlag. Berlin Heidelberg (2008) 173–214.

[10] D. Thain, T. Tannenbaum, M. Livny, Distributed Computing in Practice: the Condor Experience. J. of Concurrency and Computation: Practice and Experience. 17 (2-4) (2004) 323 – 356.

[11] Moab Adaptive Computing Suite, http://www.adaptivecomputing.com/products/moab-adaptive-computing-suite.php.

[12] D. Jackson, Q. Snell, M. Clement, Core Algorithms of the Maui Scheduler, Springer, Heidelberg, LNCS 2221 (2001) 87-102.

[13] V. Toporkov, A. Toporkova, A. Bobchenkov, D. Yemelyanov, Resource Selection Algorithms for Economic Scheduling in Distributed Systems. Procedia Computer Science. Elsevier. 4 (2011) 2267-2276.

[14] V. Toporkov, D. Yemelyanov, A. Toporkova, A. Bobchenkov, Resource Co-allocation Algorithms for Job Batch Scheduling in Dependable Distributed Computing. Dependable Computer Systems. Springer-Verlag, AICS. Berlin, Heidelberg. 97 (2011) 243-256.

[15] V.Toporkov, A. Bobchenkov, A. Toporkova, A. Tselishchev, D. Yemelyanov, Slot Selection and Co-allocation for Economic Scheduling in Distributed Computing. Proc. of the 11th Intern. Conf. on Parallel Computing Technologies. Springer-Verlag, LNCS. 6873 (2011) 368–383.