

Dependable Strategies for Job-Flows Dispatching and Scheduling in Virtual Organizations of Distributed Computing Environments

Victor Toporkov, Alexey Tselishchev, Dmitry Yemelyanov,
and Alexander Bobchenkov

Abstract. This work presents dispatching strategies based on methods of job-flow and application-level scheduling in virtual organizations of distributed computational environments with non-dedicated resources. Dependable job-flow management is implemented with the set of specific rules for resource usage. Strategies are based on economic scheduling models and diverse administration policies inside resource domains. Job management structures and economic mechanisms for load balancing in distributed environments are considered. Scheduling methods composing priority algorithms for parallel applications and job batch scheduling in distributed computing with non-dedicated resources are proposed.

1 Introduction

Distributed computational environments such as Grid have been known for significant efficiency increase in shared computational resource usage and provision of scientific and enterprise communities with solutions for complex computational tasks. However, those who are responsible for setting up Grid infrastructure and economy encounter difficulties while defining policies and strategies for efficient resource management and job scheduling. The problem of establishing an optimal or at least good strategy based on current environment condition remains actual and prominent at the moment in the domain of distributed computing.

Victor Toporkov · Dmitry Yemelyanov · Alexander Bobchenkov
National Research University “MPEI”, ul. Krasnokazarmennaya 14, Moscow,
111250 Russia
e-mail: ToporkovVV@mpei.ru,
{groddenator, yemelyanov.dmitry}@gmail.com

Alexey Tselishchev
CERN (European Organization for Nuclear Research),
CERN CH-1211 Genève 23 Switzerland
e-mail: Alexey.Tselishchev@cern.ch

Heterogeneity, changing composition, different owners of different nodes whose computing time is partially shared by users turn the organization of a distributed computational environment into an especially difficult task. Utility grid [1], multi-agent systems [2] and cloud computing [3] are types of distributed environments where usage of economic mechanisms is seen as promising. Those economic mechanisms are designed to solve tasks like resource management and scheduling of user jobs in a transparent and efficient way. Within the context of any used economic model the interests of different participants of a distributed computing environment (such as end-users or node owners) are often contradictory. Since the resources of distributed environment such as Grid are non-dedicated, it is assumed that node owners may have local job flows (their own tasks) and global job flow (which is formed by external user jobs) competing for limited computational resources of the node. Elaboration of pricing rules which are used to calculate a fee for node computing time usage and take into account user-required quality of service (QoS) is also a very serious problem [1-3]. An overview of various approaches to this problem is given in [4]. Heuristic algorithms for resource selection based on user-given utility function are described in [5]. Some resource management models offer simple search and selection of resources required by a user [6] and do not support any optimization. Others do not take into account features related to global and local job competition, the competition among users and other characteristics of distributed environments with non-dedicated computational resources [7]. A resource broker model [1-5] dynamically employs various economic policies which perform resource management which is decentralized and application-specific and have two parties: node owners and brokers representing users. Another common trend is related to virtual organizations [7-9] with central schedulers providing job-flow level scheduling and optimization. While former type of resource management is well-scalable, the simultaneous satisfaction of various application optimization criteria submitted by independent users is unreachable in essence and also can deteriorate such integral quality of service rates as total execution time of a sequence of jobs or overall resource utilization. The latter type, virtual organizations naturally restrict the scalability. However, scheduling based on uniform and controlled rules for allocation and consumption of resources makes it possible to improve the efficiency of resource usage and find a tradeoff between contradictory interests of different participants.

In this work, we propose two-level model of resource management system which is functioning within a virtual organization (VO). Resource management is implemented with a hierarchical structure consisting of a metascheduler and subordinate job schedulers that are controlled by the metascheduler and in turn interact with resource managers (e.g., with batch job processing systems). The application-level optimization begins when job-flow level optimization is finished. Such a flexible structure coupled with complex metascheduling approach enables multiaspect resource management and makes possible to control dynamic priority of job execution, resource selection and provide multicriterial optimization both on the job-flow scale and for specific job, according to its submitter requirements and optimization criteria. Hence, we may speak not only of a scheduling algorithm but rather of a scheduling strategy that is a combination of various methods of

external and local scheduling. Such a mechanism allows finer control and higher overall resource management efficiency in a distributed computing environment. *Resource* is defined as an abstract computational entity, which can be used for execution of one and only one *task*. The complex set of connected interrelated tasks form a *job*. In some applications jobs require co-scheduling and resource co-allocation on several resources [10-13]. In this case resource allocation has a number of substantial specific features caused by autonomy, heterogeneity, dynamic content changes, and node failures [6-9]. In our model jobs are submitted to the system by end-users. The proposing approach is more or less the same as used in gLite Workload Management System, where Condor is used as a scheduling module [14]. But the significant difference between the approach proposed in this work and well-known scheduling solutions for distributed environments such as the Grid [1, 3-7] is the fact that the execution strategy is formed on a basis of formalized efficiency criteria, which efficiently 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. Thus, two approaches are uniquely combined in a proposed two-tier model.

This work is organized as follows. Section 2 overviews model components and metascheduling workflow. In section 3 a strategy search is formalized. Section 4 contains simulation results. Section 5 summarizes the work and describes further research topics.

2 Basic Notions and Informal Model Components Description

Let us define basic model components presented in this work.

- VO, that defines resource co-allocation dispatching strategies, pricing policies and resource load-balancing mechanisms.
- Heterogeneous hierarchical computational environment that contains computational resources (Grid nodes, CPUs or others) with different performance indices. Each resource is considered as non-dedicated (i.e. it can have its own internal schedule and these schedules are sent to application-level schedulers upon request).
- Metascheduler, which implements resource management strategies and policies of the virtual organization.
- Application-level schedulers that analyze internal job structure and schedule single tasks.

The VO in our model of distributed computational environment includes three independent parties with their own interests.

- End-users of services provided within the VO such as computation services. End-users take steps to make resource requests to the environment, according to resource performance, time and budget estimations needed for running custom user jobs.

- VO administrators that set up resource usage policies to optimize scheduling and improve load balance. The administrators control metascheduler process running in the environment which is in fact the part of VO infrastructure software. Thus they are directly responsible for managing the parameters of higher level resource management.
- Owners of computational nodes that comprise the environment network and hardware base of the distributed computing environment. The owners offer part of their nodes computing time to VO for a fee. Computational nodes provide the only type of distributed resources used in our model.

Each computational node of the heterogeneous environment is mapped to a computational *resource line* in the metascheduler resource management routine. Several resource lines are combined into a virtual resource domain. Each resource line has two static attributes which are its performance P and its base price tag F for a computing time unit. The performance is an inherent parameter of a node and the base price tag is assigned by its owner. The dynamic characteristic of a node is represented with its local schedule which is a list of slots available for reservation. This list is sent to metascheduler by request. A slot is a continuous interval of time and is described with three parameters: its start time, its length and its fee [10-12]. The fee is calculated when the metascheduler applies its pricing policies taking in account resource type, slot length etc.

A resource request is a set of a few constraints determined by a user which correspond to the properties of the respective user job. They include:

- a) minimal performance requirement for computational nodes, P_{\min} ;
- b) maximal price tag for a single timeslot, F_{\max} ;
- c) number n of simultaneously reserved timeslots;
- d) minimal slot length;
- e) the internal structure of a job as a directed acyclic graph (DAG), where vertices represent single tasks and edges represent data dependencies [13];
- f) deadline for the job execution.

A job may require more than one timeslot if it includes several segments that can be executed in parallel way, for instance. Then the user specifies the number of reserved timeslots and minimal performance requirement that applies for them all. The whole job budget is determined by the timeslot number and the maximum price per timeslot. The minimal timeslot length requires an additional explanation. This is the minimal time estimated by the user which is required to complete job execution given the performance of the nodes meet the minimal requirement P_{\min} . Hence, the metascheduler and the user share the responsibility since the probability of being run successfully for a job equally depends on primary user estimates and overall scheduling quality.

The hierarchical model of the computational environment implies two-tier scheduling (Fig. 1). On the job-flow level the set of independent jobs is distributed between resource domains according to dispatching strategies and economic

criteria. Schedule on this level is defined by the metascheduler as a slot set for each job, which is optimal in terms of a whole job set. Application-level schedulers receive the list of resources which were meant to execute the job on and a strategy, which defines the rule used to execute tasks of a concrete job. On this level an optimal slot and specific resource are defined for each single task in a job, thus, making it possible to take internal job structure into account. On the job-flow level all end-user jobs are initially submitted into the global queue. The metascheduler can manage one or more job-flows which become sub-queues of the global queue. The mechanism of distribution of jobs between job-flows can be random or based on current load and actual efficiency of scheduling in certain job-flows. Scheduling process in each job-flow is performed by identical scheduling instance. We consider a single job-flow case.

The metascheduler works in cycles which are quanta of its process. For each cycle it has following information.

1. Information about distributed computing environment as a set of resource lines.
2. The global job queue.

What it needs then is a batch of jobs which is a ranked job list and a subset of available slots for a specific virtual resource domain and a certain timeframe which is called a scheduling interval. The length of the batch and the scheduling interval are parameterized by VO administrators.

Jobs are fetched into the batch accordingly to several variables, such as the maximum price tag, deadline, and the number of failed scheduling attempts for a job. These variables being weighted and added up determine job rank according to which it takes a position closer to head or tail of a batch.

The preparation phase ends and the actual scheduling process is executed as follows (see Fig. 1).

1. The metascheduler analyzes available slots and finds an optimal slot combination to accommodate every job in a batch using economic criteria. The budget and the deadline defined by the end-user are considered during this step. The algorithms for this step were detailed in [10-12].

2. After the domain is determined metascheduler defines the strategy for each job. For example as shown on Fig. 1, the user, who has sent the job i has the higher budget than the one who has sent the job k . The strategy for i may be expressed as “*execute as soon as possible*” while the strategy for k may be expressed as “*execute as late as possible within the defined deadline*”. These jobs are later sent to application-level schedulers and the application-level scheduling begins.

3. Application-level schedulers query internal schedules for all the resources which were selected during step 2 for each job, analyze the job DAG and form a resulting schedule for every task according to the strategy from step 2. These schedules must support interruptions and delays and should be optimal in terms of the defined criteria (i.e. cost or resource load). The criterion for the job i would

be to minimize execution cost within the defined budget, criterion for the job k would be to maximize average resource load while meeting the defined deadline. As shown on Fig. 1, jobs i and k are scheduled to be executed on the same set of resources at once.

4. Application-level schedulers are guaranteeing that there are no collisions between the tasks which were scheduled during step 3 and local tasks, which may have priority over the job-flow from step 1.

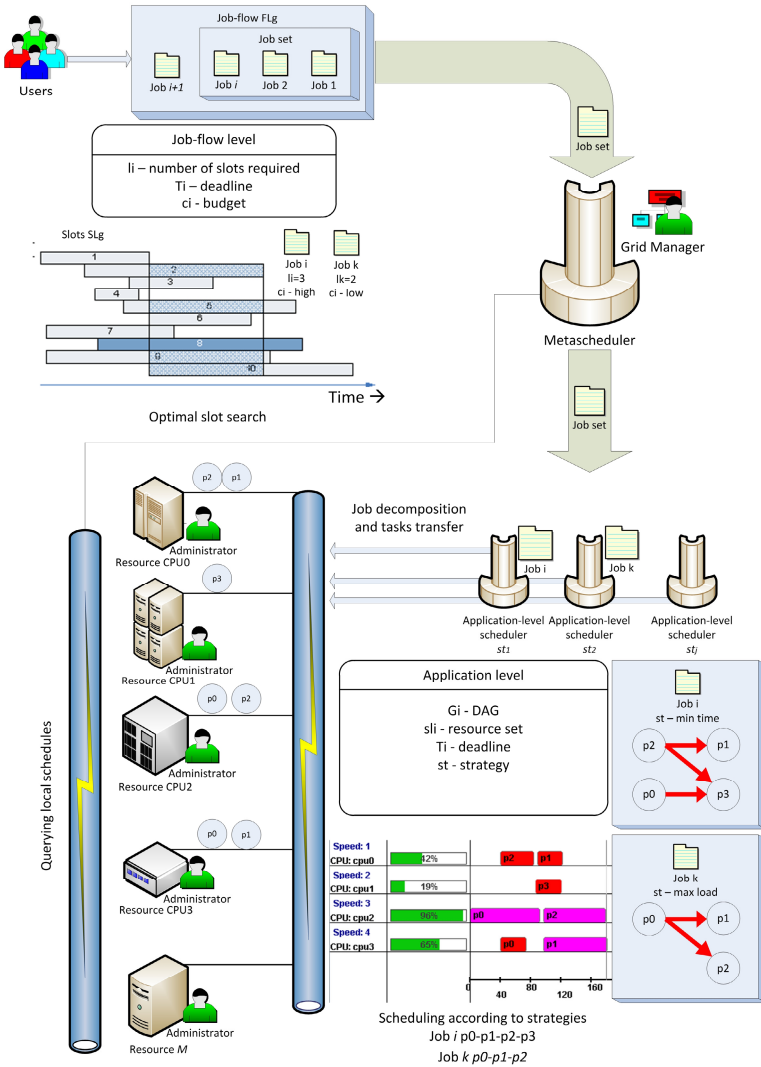


Fig. 1 Model components

3 Formalization of Scheduling

Let us note a global resource set $R_g = \{r_p, p = 1, \dots, M\}$, which includes all resources. A global job-flow is a set of jobs received by the metascheduler in time: $FL_g = \{l_i, c_i, T_i, G_i, i = 1, \dots, I\}$, where the job i is represented as l_i – the amount of resource slots required, c_i – the maximal budget end-user is ready to allocate for execution of the job, T_i – deadline, G_i – the job DAG. Metascheduler at any time moment may query each resource, receive its local schedule and build a set of slots S_{gt} – idle time intervals.

Let us introduce a set of strategies $ST = \{st_l, l = 1, \dots, L\}$, which are based on economic criteria and are defined by Grid-managers and developers. Let SL be a set of K slots suitable to execute a subset of jobs $FL_p \subseteq FL_g$. A slot set is considered as suitable for the job i if the execution is possible in terms of the resource number, the budget c_i and the deadline T_i . It is assumed that for every job there is at least one suitable slot set $sl_i \in SL, sl_i = k, k \in \{1, \dots, K\}$.

On a job-flow level for each job the metascheduler aims at finding a slot set sl_i and a strategy st_i for which the value of the function $g_i(sl_i)$, that defines whether the slot set is being effective for the job i , would be optimal [11]. The internal job structure G_i is not taken into account at this time. The mechanism to define $g_i(sl_i)$ which was developed in the previous works [10-12] is now improved. According to the resource request it is required to find a “window” with the following description: n concurrent time-slots providing resource performance rate at least P and maximal resource price not higher than F_{\max} should be reserved for a time span T_i (the resource request type was described in more detail above). The length of each slot in the window is determined by the performance rate of the node on which it is allocated. Thus as a result we have a window with a “rough right edge” (Fig. 2). In addition, the criterion of selecting the most suitable set of slots could be specified. This could be the minimum cost, the minimum runtime or, for example, the minimum power consumption criterion. The window search is performed on the list of all available system slots sorted by their start time in ascending order (this condition is necessary to examine every slot in the list and for operation of search algorithms of linear complexity [10-12]).

The scheme of a search for a window that meets the requirements and effective by the given criterion can be represented as follows.

1°. From the list of available system slots the next suitable slot s_k is extracted and examined. Slot s_k suits, if following conditions are met:

- a) resource performance rate $P(s_k) \geq P$ for slot s_k ;
- b) slot length (time span) is enough (depending on the actual performance of the slot's resource) $L(s_k) \geq T_i * P(s_k) / P$.

If conditions **a)** and **b)** are met, the slot s_k is successfully added to the window list.

2°. A current window start time is a set equal to the start time of the last added slot.

3°. Slots whose length has expired considering new window start time T_{last} are removed from the list. The expiration means that remaining slot length $L'(s_k)$, calculated like shown in **step 1°b)**, is not enough assuming the k -th slot start is equal to the last added slot start: $L'(s_k) < (T_i + (T_{last} - T(s_k)))P(s_k)/P$, where $T(s_k)$ is the slot's start time. Any combination of the remaining slots can form a window of necessary length.

4°. If the number of slots m in the current window is greater or equal to n , it is required to select n slots, effective on the specified criteria and at the same time satisfying the total cost and deadline restrictions. Suppose the window W of size n with a target criterion value equal to crW was selected. The problem of selecting efficient window consisting of n slots in the case of $m > n$ will be described below.

5°. The target criterion value crW of window W is compared with the cr' – the current best target criterion value for all previously found windows. If $crW < cr'$ (in case of a minimization problem) the window W announced as a new window-candidate and crW becomes the new best criteria value: $cr' = crW$. Go to **step 1°**.

6°. The algorithm ends after the last available slot is processed. The result of the algorithm is the window-candidate with the best target criteria value.

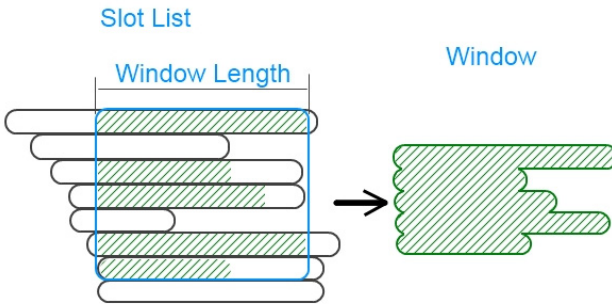


Fig. 2 Window with a “rough right edge”

The described algorithm can be compared to the algorithm of maximum/minimum value search in an array of flat values. The expanded window of size m “moves” through the ordered list of available system slots. At each step any combination of n slots inside it (in case when $n \leq m$) can form a window that meets all the requirements to run the job. The effective on the specified criteria 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 criteria value will be selected. Consider 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 case when $m = n$ the selection is trivial). The maximal budget is counted as $S = Ft_s n$, where t_s is a time span to reserve and n is the necessary number of slots. The current extended window consists of m slots s_1, s_2, \dots, s_m . The cost of using each of the slots according to their required length is: c_1, c_2, \dots, c_m . Each slot has a numeric characteristic z_i the total value of which should be minimized in the resulting window.

Then the problem could be formulated as follows:

$$a_1 z_1 + a_2 z_2 + \dots + a_m z_m \rightarrow \min, \quad a_1 c_1 + a_2 c_2 + \dots + a_m c_m \leq S,$$

$$a_1 + a_2 + \dots + a_m = n, \quad a_r \in \{0, 1\}, r = 1, \dots, m.$$

Additional restrictions can be added, for example, considering the specified value of deadline. Finding the coefficients a_1, a_2, \dots, 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 extreme value. Job-flow level scheduling ends here.

Application-level schedulers receive following input data.

- The optimal slot set sl and the description of all corresponding resources:

$$R = \{r_j, j = 1, \dots, J\} \subseteq R_g.$$

- The directed acyclic information graph $G = \{V, E\}$, where $V = \{v_i, i = 1, \dots, n\}$ is a set of vertices that correspond to job tasks, for each of those execution time estimates τ_{ij}^0 on each of resources in R are provided, E – is a set of edges that define data dependencies between tasks and data transfer time intervals.
- The dispatching strategy st , which defines the criterion for a schedule expected
- The deadline T_i or the maximal budget c_i for the job (depends on a dispatching strategy and $g_i(sl_i)$).

The schedule which is being defined on an application level is presented as follows: $Sh = \{[s_i, f_i], \alpha_i, i = 1, \dots, n\}$, where $[s_i, f_i]$ is a time frame for a task i of a job and α_i - defines the selected resource. Sh is selected in the way that the criterion function $C = f(Sh)$ achieves an optimum value. The *critical jobs method* [13] which is used to find the optimal schedule and to define f consists of three main steps.

- Forming and ranging a set of critical jobs (longest sets of connected tasks) in the DAG.
- Consecutive planning of each critical job using dynamic programming methods.
- Resolution of possible collisions.
Detailed algorithm description is presented in [13].

4 Simulation Results

The two-tier model described in the sections 2 and 3 was implemented in a simulation environment on two different and separated levels: on the job-flow level, where job-flows are optimally distributed between resource domains and on the application level, where jobs are decomposed and each task is executed in an optimal way on a selected resource.

4.1 Job-Flow Level Scheduling Simulation Results

Job-flow level metascheduling was simulated in a specially implemented and configured software that was written to test the features of the two-tier resource management.

An experiment was designed to compare the performance of our job-flow level metascheduling method with other approaches such as FCFS and backfilling. Let us remind that our scheduling method detailed in works [10] and [11] involves two stages that backfilling does not have at all, namely, slot set alternative generation and further elaboration of specific slots combination to optimize either time or cost characteristic for an entire job batch. Backfilling simply assigns “slot set” found to execute a job without an additional optimization phase. This behavior was simulated within our domain with random selection from an alternative slot, each job having one or more of them. So two modes were tested: with optimization (“OPT”) and without optimization (“NO OPT”).

The experiment was conducted as follows. Each mode was simulated in 5000 independent scheduling cycles. A job batch and environment condition was regenerated in every cycle in order to minimize other factor influence. A job batch contained 30 jobs. Slot selection was consistent throughout the experiment. If a job resource request could not be satisfied with actual resources available in the environment, then it was simply discarded.

For optimization mode as well as for no-optimization mode four optimization criteria or problems were used:

1. Maximize total budget, limit slot usage.
2. Minimize slot usage, limit total budget.
3. Minimize total budget, limit slot usage.
4. Maximize slot usage, limit slot budget.

Results presented in Table 1 apply for the problem 1. As one can see optimization mode, which is using additional optimization phase after slot set generation wins against random slot selection with about 13% gain in the problem 1 whose concern is about maximizing total slot budget thus raising total economical output per cycle and owners' profits.

Table 1 Experimental results for the problem 1: Total budget maximization with limited slot usage

| Mode | Average jobs being processed per cycle (max 30) | Average total slot cost per cycle, <i>cost units</i> | Average total slot usage per cycle, <i>time units</i> | Average slot usage limit per cycle, <i>time units</i> |
|--------|---|--|---|---|
| OPT | 20.0 | 11945.98 | 421.22 | 471.14 |
| NO OPT | 20.0 | 10588.53 | 459.36 | 471.85 |

Comparable results were obtained for other problems which are summarized in Table 2. Optimized values are outlined in light grey.

Table 2 Experimental results for the problems 2-4

| Mode | Average jobs being processed per cycle (max 30) | Average total budget (slot cost) per cycle, <i>cost units</i> | Average total slot usage per cycle, <i>time units</i> | GAIN, % |
|--|---|---|---|---------|
| Problem 1: Maximize total budget, limit slot usage | | | | |
| OPT | 20.0 | 11945.9 | 421.2 | +12.8 |
| NO OPT | 20.0 | 10588.5 | 459.4 | |
| Problem 2: Minimize slot usage, limit total budget | | | | |
| OPT | 12.4 | 7980.4 | 300.9 | +10.6 |
| NO OPT | 12.4 | 7830.9 | 332.8 | |
| Problem 3: Minimize total budget, limit slot usage | | | | |
| OPT | 15.1 | 9242.4 | 410.057 | +6.2 |
| NO OPT | 15.3 | 9813.9 | 406.612 | |
| Problem 4: Maximize slot usage, limit total budget | | | | |
| OPT | 15.28 | 9870.8 | 416.835 | +3.0 |
| NO OPT | 15.4 | 9718.1 | 404.8 | |

These results are showing the advantage of the metascheduling on the job-flow level. The next section describes the experiments on the application level.

4.2 Application Level Scheduling Simulation Results

The experiment results presented in Table 3 shows the advantage of the critical jobs method usage in a two-tier scheduling model compared to consecutive application-level scheduling. Here $k=0.75$ means that each job is sent to be scheduled after 75% of the time allocated for the previous one: while the scheduling cost for a job is more or less the same, 1000 jobs are planned 25% faster.

Consider another experiment: while changing the length of the scheduling interval, we will estimate the proportion of successfully distributed jobs. The length of the scheduling interval is equal to $L = l * h, h = 1.0, \dots, 2.6$, with step 0.2, where l is the length of the longest critical path of tasks in the job and h is a distribution interval magnification factor. There were carried 200 experiments for each h (bold points on Fig. 3). Analysis of the Fig. 3 shows that increasing the scheduling interval (relatively to the execution time of the longest critical path on the nodes with the highest performance) is accompanied by a significant increase in the number of successfully distributed jobs. The detailed study of this dependence can give a priori estimates of an individual job successful distribution probability.

Table 3 Two-tier model vs consecutive application-level scheduling

| Parameter | Application-level scheduling | Two-tier model ($k=0.75$) |
|----------------------|------------------------------|-----------------------------|
| Jobs number | 1000 | 1000 |
| Execution time | 531089 time units | 399465 time units |
| Optimal schedules | 687 | 703 |
| Mean collision count | 3.85 | 4.41 |
| Mean load (forecast) | 0.1843 | 0.1836 |
| Mean load (fact) | 0.1841 | 0.1830 |
| Mean job cost | 14.51 units | 14.47 units |

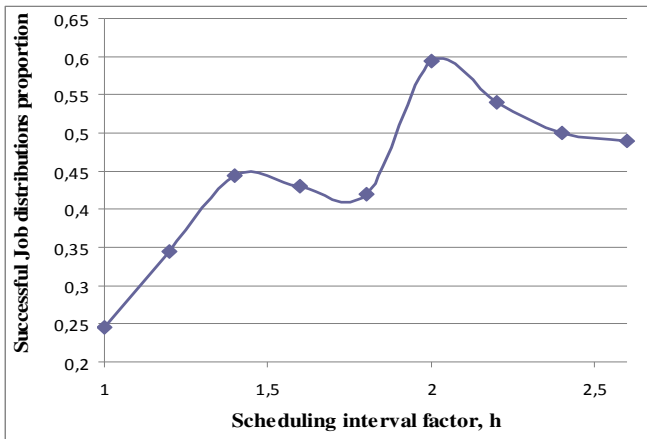
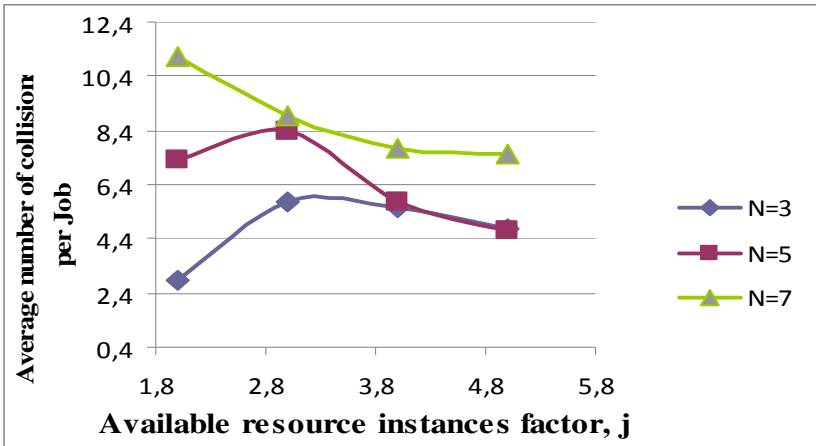


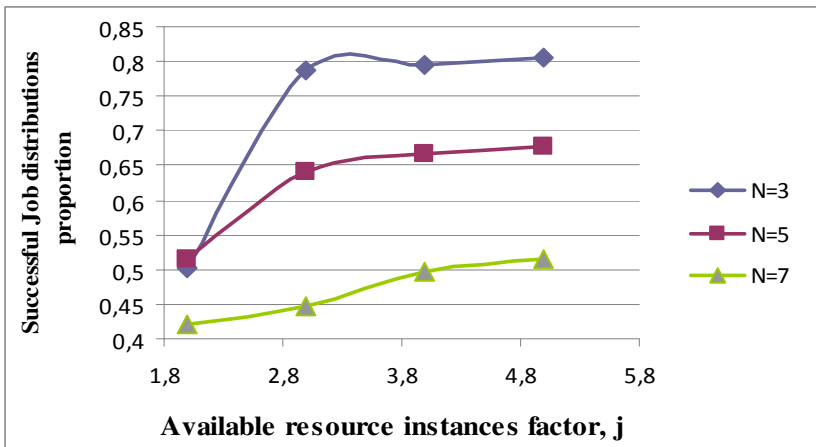
Fig. 3 Dependence of the proportion of the successful job distributions on the length of the distribution interval

In the next experiment we will consider the dependence of successful distributions number and the number of collisions per experiment on the level of resource instances availability. The experiments were performed in conditions of limited resources using the specific instances of the resources. The number of resources J in each experiment was determined as $J = j * N$, where j – factor (x-axis) and N – number of tiers in the graph. Fig. 4 shows results of the experiments with different j values and $N = 3, 5, 7$.

The obtained dependencies (Fig. 4) suggest that the collisions number depends on the resources availability. The lower the number of resource instances and the greater the number of tiers in the graph – the more collisions occurred during the



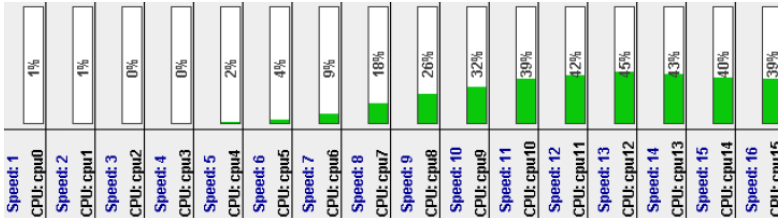
(a)



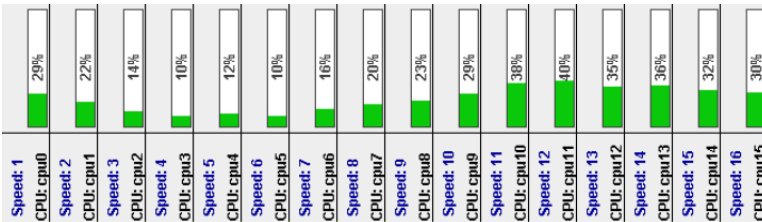
(b)

Fig. 4 Simulation results: resource dependencies of collisions number (a) and successful job distribution proportion (b)

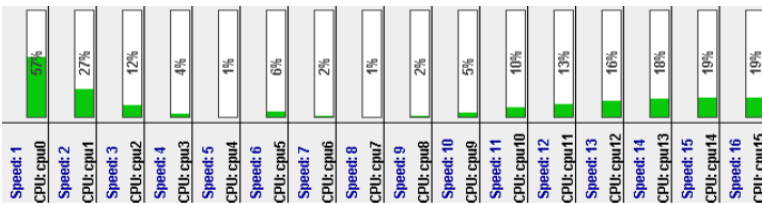
scheduling. At the same time the number of resource instances affects the successful distribution probability. With a value of $j > 4$ (that is, when the number of available resource instances is more than 4 times greater than the number of tiers in the graph) all cases provide the maximum value of successful distribution probability. These results are subject of future research of refined strategies on a job-flow level.



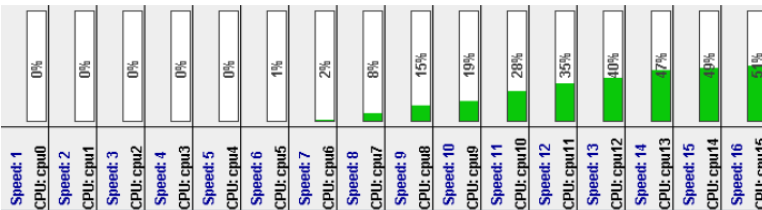
(a)



(b)



(c)



(d)

Fig. 5 Resource utilization level balancing: utilization maximization with $h = 1.66$ (a) and $h = 1.2$ (b), utilization minimization and distribution cost maximization (c), distribution cost minimization (d)

The next series of experiments aimed at identifying the priorities of selecting certain resource instances with different optimization criteria and various restrictions. Figures 5 (a-d) show resource utilization levels in the following problems: resource load balancing, distribution cost minimization and maximization. The scheduling interval is defined as $L = l * h$, $h = 1.66$ and 1.2 , where l is execution time of the longest critical path and h is a scheduling interval factor. Processors with greater number have relatively lower cost and performance level. To maximize average resource utilization the priority is given to processors with relatively low performance (Fig. 5 (a)). In case of a shorter scheduling interval ($h = 1.2$) there is need to use resources with higher performance (Fig 5 (b)). During the resource utilization minimization and distribution cost maximization the priority is given to the nodes with higher performance and usage cost (Fig. 5 (c)). During the distribution cost minimization the priority is given to processor with low performance level and correspondingly low cost. These experiments show how strategies defined on a job-flow level are implemented on an application level, how flexible the strategies can be and how can resource load be controlled by the metascheduler.

5 Conclusions and Future Work

In this work, we address the problem of independent job-flow scheduling in heterogeneous environment with non-dedicated resources.

Each job consists of a number of interrelated tasks with data dependencies. Using the combination of existing methods with a number of original algorithms the resulting schedules are computed. These schedules meet the defined deadlines and budget expectations, provide optimal load-balance for all the resources and follows virtual organization's strategies, thus, allowing to achieve unprecedented quality of service and economic competitiveness for distributed systems such as Grid. The experiments which were conducted are showing the efficiency of methods developed for both job-flow and application level scheduling. The model proposed is showing the way these methods and advantages can be converged in one place making it possible to achieve the main goal.

Future research will include the simulation of connected job-flow and application levels and experiments on real Grid-jobs in order to get finer view on advantages of the approach proposed.

Acknowledgements. This work was partially supported by the Council on Grants of the President of the Russian Federation for State Support of Leading Scientific Schools (SS-316.2012.9), the Russian Foundation for Basic Research (grant no. 12-07-00042), and by the Federal Target Program "Research and scientific-pedagogical cadres of innovative Russia" (State contracts 16.740.11.0038 and 16.740.11.0516).

References

- [1] Garg, S.K., Buyya, R., Siegel, H.J.: Scheduling parallel applications on utility Grids: time and cost trade-off management. In: Proc. of ACSC 2009, Wellington, New Zealand, pp. 151–159 (2009)

- [2] Tesauro, G., Bredin, J.L.: Strategic sequential bidding in auctions using dynamic programming. In: Proc of the First International Joint Conference on Autonomous Agents and Multiagent Systems: part 2, pp. 591–598. ACM, New York (2002)
- [3] Garg, S.K., Yeo, C.S., Anandasivam, A., Buyya, R.: Environment-conscious scheduling of HPC applications on distributed cloud-oriented data centers. *J. of Parallel and Distributed Computing* 71(6), 732–749 (2011)
- [4] Buyya, R., Abramson, D., Giddy, J.: Economic models for resource management and scheduling in Grid computing. *J. of Concurrency and Computation: Practice and Experience* 14(5), 1507–1542 (2002)
- [5] Ernemann, C., Hamscher, V., Yahyapour, R.: Economic Scheduling in Grid Computing. In: Feitelson, D.G., Rudolph, L., Schwiegelshohn, U. (eds.) JSSPP 2002. LNCS, vol. 2537, pp. 128–152. Springer, Heidelberg (2002)
- [6] Voevodin, V.: The Solution of Large Problems in Distributed Computational Media. *Automation and Remote Control. Pleiades Publishing, Inc.* 68(5), 773–786 (2007)
- [7] Kurowski, K., Nabrzyski, J., Oleksiak, A., et al.: Multicriteria aspects of Grid resource management. In: Nabrzyski, J., Schopf, J.M., Weglarz, J. (eds.) *Grid Resource Management. State of the Art and Future Trends*, pp. 271–293. Kluwer Acad. Publ. (2003)
- [8] Toporkov, V.: Application-Level and Job-Flow Scheduling: An Approach for Achieving Quality of Service in Distributed Computing. In: Malyskin, V. (ed.) PaCT 2009. LNCS, vol. 5698, pp. 350–359. Springer, Heidelberg (2009)
- [9] Toporkov, V.V.: Job and application-level scheduling in distributed computing. *Ubiquitous Comput. Commun. J.* 4, 559–570 (2009)
- [10] Toporkov, V., Toporkova, A., Bobchenkov, A., Yemelyanov, D.: Resource selection algorithms for economic scheduling in distributed systems. *Procedia Computer Science* 4, 2267–2276 (2011)
- [11] Toporkov, V., Yemelyanov, D., Toporkova, A., Bobchenkov, A.: Resource Co-allocation Algorithms for Job Batch Scheduling in Dependable Distributed Computing. In: Zamojski, W., Kacprzyk, J., Mazurkiewicz, J., Sugier, J., Walkowiak, T. (eds.) *Dependable Computer Systems. AISC*, vol. 97, pp. 243–256. Springer, Heidelberg (2011)
- [12] Toporkov, V., Bobchenkov, A., Toporkova, A., Tselishchev, A., Yemelyanov, D.: Slot Selection and Co-allocation for Economic Scheduling in Distributed Computing. In: Malyskin, V. (ed.) PaCT 2011. LNCS, vol. 6873, pp. 368–383. Springer, Heidelberg (2011)
- [13] Toporkov, V.V., Tselishchev, A.S.: Safety scheduling strategies in distributed computing. *Intern. J. of Critical Computer-Based Systems* 1(1/2/3), 41–58 (2010)
- [14] Cecchi, M., Capannini, F., Dorigo, A., et al.: The gLite Workload Management System. *Journal of Physics: Conference Series* 219(6), 062039 (2010)