

# Markets as Global Scheduling Mechanisms: The Current State

Junko Nakai

AMTI at NASA Ames Research Center, Mail Stop 258-6, Moffett Field, CA  
94305-1000  
nakai@nas.nasa.gov

**Abstract.** Viewing the computers as “suppliers” and the users as “consumers” of computing services, markets for computing services/resources have been examined as one of the most promising mechanisms for global scheduling. We first establish how economics can contribute to scheduling. We further define the criterion for a scheme to qualify as an application of economics. Many studies to date have claimed to have applied economics to scheduling. If their scheduling mechanisms do not utilize economics, contrary to their claims, their favorable results do not contribute to the assertion that markets provide the best framework for global scheduling. For any of the schemes examined, we could not reach the conclusion that it makes full use of economics.

## 1 Scheduling and Economics

A distributed computing system brings about opportunities for enhancing the efficiency in computing, and hence, for increasing the value of existing facilities for computation. One of the new administrative features necessitated for that purpose is global scheduling. The core problem we face in scheduling computing jobs is exceedingly similar to that in the economy. If there were an unlimited amount of computing resources, jobs would be executed as soon as they are submitted, eliminating the need for scheduling. As there is a limit to the availability of resources, scheduling of computing jobs becomes a problem of allocating limited resources. We may envision a situation where the users of computers are required to give up the limited resources they are endowed with in exchange for access to the computing resources. Further, if the endowment could be used for more than one item or occasion, which result in outcomes differing in importance, then the allocation of computing resources is precisely the economic problem in its most fundamental form. Modern economics is recognized as the science that studies human behavior as a relationship between scarce means which have alternative uses, as asserted by Robbins [17].<sup>1</sup>

---

<sup>1</sup> The most basic activity in an economy is an exchange of goods and services. Whenever there is an exchange of goods or services, a price is established, which is simply the rate of exchange of goods or services involved. Thus, when we refer to one of the three, an economy, an exchange, and a price, the other two necessarily exist. Markets are forums for exchanges of goods and services, where exchanges are voluntarily

When the so-called pricing of computing resources started to attract increasing attention in the late 1960s, it was quickly acknowledged that pricing can be seen as a kind of scheduling mechanism and that both concern allocation of resources [16].<sup>2</sup> The advent of networks of computers did not alienate scheduling from economics, and scheduling for distributed computer systems has also been recognized as an activity under resource management [2, 12]. Markets for computing resources have been embraced by many as the key components in the operation of distributed systems [1, 3, 4, 5, 10, 11, 14, 15, 18, 19, 20, 21, 22, 23].

### 1.1 The Objective of an Economy

An additional insight to the economic problem, as defined by Robbins, was provided by Hayek [6]. He made it clear that the problem was how to allocate resources so that they are used in the best way possible. Inasmuch as the best uses are known only to individuals, the economic problem becomes “a problem of the utilization of knowledge which is not given to anyone in its totality” [6]. By pursuing the “best” use of each resource, our attempt to allocate resources becomes one of maximizing private values associated with resource use. Certainly, this is also the problem posed to the group of users who are to share a set of computing facilities. The computing capabilities should be shared, and for that reason, information on the requirements for each job, including what resources and the value of its successful execution, is required. Unfortunately, this information is usually available only to the user. In fact, utilization of users’ private information has been recognized as one of the advantages of pricing over other types of scheduling [16].

We now discuss the importance of scarcity of resources, including budget, and that of the existence of alternative uses of resources in the formation of value itself. A variety of needs, which cannot all be fulfilled and each of which leads to an outcome that is different in importance, is what brings the economic problem into its existence. If resources are obtainable whenever desired and/or their uses result in outcomes identical in value, there would be no question of allocating the resources for their best use or to maximize value from use.

We distinguish two types of scarcity. It may take the straightforward form of a finite limit to the amount available, an amount smaller than is required to fulfill all needs, in which case the multiplicity of possible uses of the resource, with outcomes that are distinguishable in importance, leads to the economic problem. If a limited amount of a resource is available during a particular time period (which cannot be used any other time) and that resource is designated for only one use, the economic problem does not exist because there is a unique possibility: allocation of the entire resource for that single use. Scarcity may

---

initiated. Markets do not prevail in all types of economies, and economies are not synonymous with markets.

<sup>2</sup> Pricing stepped into the limelight because the most common default scheduling mechanism, first-come, first-served, was often combined with additional rules, indicating that it alone was not meeting needs [16].

take another form: resources having dual effects, each of which is associated with a positive or a negative value, depending on the amount allocated. The second type of scarcity also serves the function provided by alternative uses, which vary in importance of attainable outcomes, in the first kind. Both the existence of alternative uses with different, resultant utility levels (in the first kind of scarcity), and the possible negative effects of resource use (in the second type of scarcity) make it necessary for resource users to evaluate and compare the values of various allocations.

## 1.2 Application of the Economics to Scheduling

We have laid out above the problem an economy faces, and clarified the motivation behind economic activities. The existence of a motivation to meet an objective is a necessary condition for the attainment of the objective. In the following sections, we analyze various pricing schemes for computing resources, with this necessary condition as a yardstick for determining whether a scheme is an application of economics. When the participants in the scheme maximize utility (or equivalently, value), given constraints, by choosing among needs that cannot be simultaneously fulfilled and are different in achievable utility, the economic motivation exists, and therefore, the necessary condition is satisfied for calling the scheme an application of economics to scheduling. If a computing-service provider is directly involved in the process of determining service allocation (i.e., the provider interacts with users in the process), the criterion must be satisfied not only by the users but also by the provider.

## 2 Markets and Computer Networks: from early 1980s to present

A computer network consists not only of multiple users, but also of multiple service providers, appearing much more complex [19] and closer to a market as we know from everyday life [1, 5, 10, 11, 15, 21]. The seeming ease with which markets allocate resources, albeit its complexity [1, 5, 11, 15, 19], and the solid existence of theoretical microeconomics (particularly, the general equilibrium theory) [5, 11, 15, 19], have captured the imagination of researchers in the field of networks. We examine below the representative market models for distributed computing systems (the Contract Net Protocol [18], the Enterprise [14], a model by Kurose and Simha [11], Agoric System [15], Spawn [10, 21], WALRAS [3, 22, 23], Mariposa [19, 20], and the Grid Architecture for Computational Economy [1]),<sup>3</sup> focusing on the ones that make an extensive and direct use of the general equilibrium theory: the model by Kurose and Simha, and WALRAS.

---

<sup>3</sup> In citing such models, those for database systems that assume light traffic, are often included, e.g., Mariposa.

## 2.1 Model by Kurose and Simha

The paper by Kurose and Simha [11] is one of the oft-cited papers on distributed computer systems with relation to economics. We briefly discuss the general equilibrium theory, which their model draws on, with attention to the so-called tâtonnement process and the problematic aspects of the process. While the second type of scarcity existed for the nodes in the model, they were not utility maximizers, disqualifying the model as an application of economics. The utility of the resource allocator, called auctioneer, was undefined.

**General Equilibrium Theory and Resource-Directed Approach.** Kurose and Simha implemented one of what they call the two basic microeconomic approaches, the price-oriented and the resource-oriented approaches, which are better known as a tâtonnement process (for reaching an equilibrium) with pure price adjustment and that with pure quantity adjustment, respectively. Hence, the model applies the general equilibrium theory in economics, a theory that concerns price and quantity determination in equilibrium initiated by Léon Walras in the late 19th century. The economy under consideration in the general equilibrium theory is one in which the number of agents in the system is large enough so that any one of them cannot affect prices by acting alone, and agents do not collaborate.<sup>4</sup>

The tâtonnement process with pure quantity adjustment in a pure exchange economy proceeds as follows: The auctioneer informs the utility-maximizing user-agents of their entitled quantities of resources and the agents report back the marginal utilities at those quantities. The auctioneer changes the allocation of resources so that more inputs are allocated to the agents with higher marginal utilities. The process continues until an equilibrium in price and quantity is attained. This is the approach preferred and adopted by the authors, because all interim allocations are feasible, unlike the price-oriented approach. We note that the process does not guarantee convergence to a unique equilibrium, even if one exists, without further restrictions on the economy.

Another attractive feature of the resource-directed approach was reported: “When analytic formulas are used to compute performance, successive iterations of the algorithm result in resource allocations of strictly increasing systemwide utility.”<sup>5</sup> Putting aside the issue of incentive compatibility, we may say that optimization by local agents led to an optimal solution for the entire system

---

<sup>4</sup> In other words, all participants in the economy, producers and consumers, are price-takers.

<sup>5</sup> Heal [7] concluded that in the tâtonnement process with pure quantity adjustment more information exchange would be required than in the process with pure price adjustment. Hurwicz [9] pointed out that the total information would be of higher dimension in the process with pure quantity adjustment compared to the process with pure price adjustment, but also added that whether the difference was significant was “somewhat controversial.”

because the global utility was set equal to the sum of utilities of local agents.<sup>6</sup> The iterations in resource allocation conform with the economic motivation (or, are economically feasible) in the adopted framework, only under the condition that honest reporting of utility levels is ensured.

The two advantages described above, feasibility (in a discrete process) and monotonicity, were labeled two desirable properties of tâtonnement [13]/gradient [9]-based processes for reaching an equilibrium, by Malinvaud.<sup>7</sup> Kurose and Simha proposed algorithms which were based on a tâtonnement process with pure quantity adjustment. Theoretical investigation of their most basic algorithm by Heal [7] had shown that it indeed exhibits both properties. For the process to function, however, one condition has to be met: the central authority's knowledge of an initial allocation that is feasible, which is the cost of obtaining the desirable properties according to Hurwicz [9].

**Optimal File Allocation.** The distributed system chosen for investigation was a network of nodes, which were assumed capable of communicating with any other in the network. The problem was how to allocate files optimally.<sup>8</sup> The cost of communication to each node was defined to be the average delay in the transmission of messages. The cost of access delay was defined to be the expected time in access delay. In turn, the sum of the cost of communication and the cost of access delay was called the expected cost of access to the file source at a node. Therefore, allocating all files at one node may reduce the cost of communication, but only by increasing the cost of access delay, because that node must handle all inquiries in the system. The expected cost of access to the entire network was the sum of the expected cost at each node. The optimal file allocation was taken to be the allocation that minimizes the expected cost of access for the whole network.<sup>9</sup>

The performance of three decentralized algorithms was examined. They all employed gradient processes; each node computed the first derivative of the utility function (i.e., marginal utility) and/or the second derivative, evaluated at a specific point, and sent that information to the central node (or alternatively, to all other nodes). The paper concluded that all algorithms had the following desirable properties: feasibility in all iterations, strict monotonicity, and fast convergence.

---

<sup>6</sup> Consequently, the global objective function was a function only of local objective functions increasing in all of its arguments, and unresponsive to the names of the local agents; local optimization coincided with global optimization.

<sup>7</sup> He was concerned about the possibility of slow or "disorganized" convergence to an equilibrium, which implied that the existence of an equilibrium and convergence to one through a tâtonnement process, by themselves, do not guarantee practicality of the process as one in economic planning.

<sup>8</sup> Each node had a local look-up table, which provided information on the file fragment locations so that a request not met locally could be sent to an appropriate node.

<sup>9</sup> In order to cast the problem as a utility maximization problem, the utility was set equal to the negative of the expected cost.

**Problems in General Equilibrium Theory** One of the problems in the general equilibrium theory is that the tâtonnement process cannot do without an auctioneer; decentralized economic planning does not truly qualify as a decentralized system with features such as lack of a single point of failure, as envisaged by many of the researchers in the field of global scheduling. Heal's model [7], on which the model by Kurose and Simha is based, was decentralized only in the sense that information was collected from the local agents. The mechanisms employed by Kurose and Simha for adjustment of the system towards an equilibrium are informationally decentralized [9], but that is not equivalent to decentralized decision-making.<sup>10</sup> Indeed, the processes do not concern local decision-making, just as Heal's model does not.

Another problem in the general equilibrium theory is also carried over to the model by Heal, as well as to that by Kurose and Simha: incentive incompatibility.<sup>11</sup> Note that individual nodes could have increased the final utility by falsely reporting the levels of marginal utility that were above the actual levels. However, the local agents in the investigated network acted so as to fulfill the global goal, by forgoing the opportunity to increase their own utilities. Moreover, the utility of auctioneers is always ignored in the general equilibrium theory, and so it is in the models by Heal, Kurose and Simha.

The above discussion also serves as an analysis of whether the resource allocation scheme was driven by economic considerations of the agents in the system. The model relied on nodes' balancing the benefits and the costs of owning a file fragment; the second type of scarcity (and hence, the existence of alternative uses in the broad sense) was present. The fact that the nodes reported their marginal utilities honestly to the central node (or, to all other nodes), in face of feasible cheating and attainment of higher utility, together with the fact that honesty did not factor into the utility defined, indicate that the nodes were not utility-maximizing agents; the first part of the necessary condition to be an economic application was not met. The utility of the auctioneer, an active participant in the resource allocation process, was left undefined. The study neither validates nor invalidates the appropriateness of creating markets for computing resources in distributed systems.

## 2.2 WALRAS

WALRAS, like the model by Kurose and Simha, is an application of the general equilibrium theory to scheduling [22]. It implemented the tâtonnement process with pure price adjustment. WALRAS is incentive incompatible, as any model based on the general equilibrium theory with a finite number of agents would

<sup>10</sup> An informationally decentralized process is one which has informational requirements that are no greater than those for a perfectly competitive process [8].

<sup>11</sup> Although reporting the derivatives of the pertinent functions to the central authority is one of the standard elements in the tâtonnement processes, its incentive compatibility has been established only when the number of traders is infinite in a pure exchange economy, if no forced, initial redistribution of endowments is allowed [9].

be (if there exists neither production nor the possibility of redistributing initial endowments); it is unsuitable to be called an economic application.

**General Equilibrium Theory and WALRAS.** In the standard tâtonnement process with pure price adjustment, there is an auctioneer who informs agents of the prices, and the agents report back the amounts of goods they demand (or more precisely, demand for goods over and above the amount endowed) at those prices. Such reports are called bids. The auctioneer calculates the new prices, according to the predetermined rules and the amounts of demand reported, and the process repeats until the prices no longer need to be adjusted.

WALRAS differed from the standard tâtonnement process in that demand functions were reported by the agents (instead of a point on a demand function) and that the auctioneer dealt with each good separately [3, 22]. Moreover, not all agents reported their demands for all goods in each time period [3]. Random draws, which were independent across time and agents, determined which bids were submitted. For the unselected combinations of agents and goods, the bids from the previous period were used. The advantage of their asynchronous bidding was small price oscillations [3, 23]. Cheng and Wellman [3] favored their approach over other processes for attaining an equilibrium in the general equilibrium theory, since they saw fewer opportunities for strategic interactions and no trade took place until an equilibrium was achieved (no resource was allocated based on intermediate results, which were by definition not global optima and may have been irreversible if implemented). The utility function of the resource users,  $u(x)$ , was of constant elasticity of substitution:  $u(x) = \left( \sum_{j=1}^k \alpha_j (x^j)^\rho \right)^{1/\rho}$ , where  $\alpha_j$ 's were randomly generated coefficients from a uniform distribution,  $x^j$  was the amount of good  $j$ , and  $\rho$  was fixed at 0.5 (for the main simulation). Therefore, the resulting excess-demand functions, for each agent and for the entire economy, had the property of gross substitutability [3]. The existence of an equilibrium was guaranteed by the preferences implied from the utility functions (which were continuous, strictly convex, and locally nonsatiated) and non-negative total excess-demand (as Cheng and Wellman implicitly demonstrated with experiments).<sup>12</sup> Moreover, gross substitutability ensured the uniqueness of equilibrium and convergence to that equilibrium point on any price path. While the adaptive learning behavior of the auctioneers justified the rules of WALRAS, users remained simple price-takers who reported their excess-demand functions honestly. The users would not have reported the true demand functions if they acted so as to maximize utilities, as is evidenced by the utility function shown above.

**Convergence and Other Problems.** We concentrate here on the most comprehensive results given in “The WALRAS Algorithm” [3]. An examination of 100 randomly generated economies, with five or seven agents of utility as described above (where  $j$  is equal to 5), showed that the median behavior of the

<sup>12</sup> See Figure 1 in “The WALRAS Algorithm” [3].

system was a rapid convergence to the equilibrium at the beginning, and leveling off at a small, positive amount of total excess-demand after 150 iterations.<sup>13</sup> When values other than 0.5 for the substitution coefficient were adopted, convergence was not seen even after 5,000 iterations in some cases.

The feasibility of the proposed algorithm and singularity of equilibrium are obtained only under restricted circumstances [3]; we need restrictions on the agents' utility functions. Gross substitutability of the aggregate excess-demand function is a sufficient condition, and the authors reported that they could not find a class of utility functions that are not grossly substitutable and yet converge to an equilibrium. Whatever the necessary conditions may be for convergence, the utility functions employed must represent the preferences of users. WALRAS explored the case where every agent had an identical utility function. How it would serve a system of users with various utility functions and how it would compare with other scheduling mechanisms are yet to be seen.

### 2.3 Other Models

The information exchange process in the Contract Net Protocol [18] does not appear too different from that in a non-distributed computing system, if evaluated according to the characterization provided by the author of the protocol. Many of the details necessary for implementation were left unspecified, including agents' utility; we are unable to conclude that the protocol is an economic application.

The Enterprise system [14] is a fleshed-out version of the Contract Net Protocol, which connected personal workstations, using a local area network. Although favorable simulation results were reported, we cannot attribute them wholly to the scheme's general features. Neither can we conclude that the scheme is an application of economics, because the utility of the service providers in the scheme was undefined. There was no mention of user budget, without which we cannot evaluate whether the scarcity and alternative-use condition for users was satisfied.

Spawn [10, 21] is a resource allocation mechanism for a network of heterogeneous workstations whose agents are sellers (i.e., owners of workstations, who are not using them at any given moment) and buyers of CPU time. The special feature of Spawn is its spawning process or dividing a task into subtasks. The agents' utilities (that of both buyers and sellers of computing resources) were not defined; it was impossible to confirm that the scheme satisfied the necessary condition to be an economic application.

Mariposa is a distributed database and storage system for non-uniform, multi-administrator, wide area networks [19, 20]. Some of the advantages of markets, which are claimed also as those of Mariposa, do not hold unconditionally, and the practicality of the objectives chosen for the artificial agents is

<sup>13</sup> Equilibria reached through tâtonnement processes in a framework as the one adopted by WALRAS are Pareto-optimal, which are often interpreted as desirable allocations [23].



in question. Some, but not all, of the agents' objectives (and thus, utilities of some agents) were proposed. We could not conclude that the scheme satisfied the necessary condition for being an economic application.

An agoric system is an intellectual exploration as to what economics may be able to do for distributed computing systems. Without more information than has been provided in the paper [15], we cannot draw a conclusion as to whether the system is an application of economics.

The Grid Architecture for Computational Economy [1] aims at incorporating an economic model into a grid system with existing middleware, such as Globus and Legion. Although utilities for service providers and users were suggested, there was an implication that no alternative uses for user budget existed; the scheme does not qualify as an economic application.

### 3 Future Direction

We could not conclude that the criterion for an economic application was satisfied by any of the models examined. Moreover, there was no comparison of their performances with those using other scheduling mechanisms. Hence, no support was provided for establishing that economics is a necessary component for superior performance of distributed computing systems.

Many of the studies examined justified their use of economics based on the desirable results given by the general equilibrium theory, which is more narrowly focused than the whole discipline of economics. How relevant is the general equilibrium theory to scheduling? More generally, does the theory capture what drives the economy to behave well as a system? Conversely, are markets capable of producing the favorable outcomes as the general equilibrium theory implies (and subsequently, asserted by the architects of market-based scheduling)? We believe that many of these questions can be answered through a careful reading of the general equilibrium theory and its related fields.

### References

1. Buyya, R., Abramson, D., Giddy, J.: A Case for Economy Grid Architecture for Service Oriented Grid Computing. Presented at the 10th Heterogeneous Computing Workshop, San Francisco, April 23, 2001
2. Casavant, T. L., Kuhl, J. G.: A Taxonomy of Scheduling in General-Purpose Distributed Computing Systems. *IEEE Trans. Soft. Eng.* **14** (1988) 141-154
3. Cheng, J. Q., Wellman, M. P.: The WALRAS Algorithm: A Convergent Distributed Implementation of General Equilibrium Outcomes. *Computational Econ.* **12** (1998) 1-24
4. Clearwater, S. H.: Why Market-Based Control? In: Clearwater, S. H. (ed.): *Market-Based Control* World Scientific Publishing, Singapore (1996)
5. Ferguson, D. F., Nikolaou, C., Sairamesh, J. Yemini, Y.: Economic Models for Allocating Resources in Computer Systems. In: Clearwater, S. H. (ed.): *Market-Based Control*. World Scientific Publishing, Singapore (1996)

6. Hayek, F. A.: The Use of Knowledge in Society. *Amer. Econ. Rev.* **35** (1945) 519-530
7. Heal, G.: Planning without Prices. *Rev. Econ. Stud.* **36** (1960) 347-362
8. Hurwicz, L.: On Informationally Decentralized Systems. In: McGuire, C. B. Radner, R. (eds.): *Decision and Organization* North-Holland Publishing, New York (1972)
9. Hurwicz, L.: The Design Mechanisms for Resource Allocation. *Amer. Econ. Rev.* **63** (1973) 1-30
10. Huberman, B. A., Hogg, T.: Distributed Computation as an Economic System. *J. Econ. Persp.* **9** (1995) 141-152
11. Kurose, J. F., Simha, R.: A Microeconomic Approach to Optimal Resource Allocation in Distributed Computer Systems. *IEEE Trans. Comp.* **38** (1989) 705-717
12. MacKie-Mason, J. K., Varian, H. R.: Pricing the Internet. In: Kahin, B., Keller, J. (eds.): *Public Access to the Internet*. MIT Press, Cambridge, Massachusetts (1995)
13. Malinvaud, E.: Decentralized Procedures for Planning. In: Malinvaud, E., Bacharach, M. O. L. (eds.): *Activity Analysis in the Theory of Growth and Planning*. Macmillan, London (1967)
14. Malone, T. W., Fikes, R. E., Grant, K. R., Howard, M. T.: Enterprise: A Market-like Task Scheduler for Distributed Computing Environments. I; Huberman, B. A. (ed.): *The Ecology of Computation*. Elsevier Science Publishers, North-Holland (1988)
15. Miller, M. S., Drexler, K. E.: Markets and Computation: Agoric Open Systems. <http://www.agorics.com/agoricpapers.html>, 7 July 2000
16. Nielsen, N. R.: The Allocation of Computer Resources—Is Pricing the Answer? *Comm. ACM.* **13** (1970) 467-474
17. Robbins, L. C.: *An Essay on the Nature and Significance of Economic Science*. Macmillan, London (1984, originally published in 1932)
18. Smith, R. G.: The Contract Net Protocol: High-Level Communication and Control in a Distributed Problem Solver. *IEEE Trans. Comp.* **c-29** (1980) 1104-1113
19. Stonebraker, M., Devine, R., Kornacker, M., Litwin, W., Pfeffer, A., Sah, A., Staelin, C.: An Economic Paradigm for Query Processing and Data Migration in Mariposa. <http://sunsite.berkeley.edu/Dienst/UI/2.0/Describe/ncstrl-ucb/S2K-94-49>
20. Stonebraker, M., Aoki, P. M., Litwin, W., Pfeffer, A., Sah, A., Staelin, C., Yu, A.: Mariposa: a Wide-Area Distributed Database System *VLDB J.* **5** (1996) 48-63
21. Waldspurger, C. A., Hogg, T., Huberman, B. A., Kephart, J. O., Stornetta, S.: Spawn: A Distributed Computational Economy. *IEEE Trans. Soft. Eng.* **18** (1992) 103-117
22. Wellman, M. P.: Market-Oriented Programming Environment and its Application to Distributed Multicommodity Flow Problems. *J. Art. Intel. Res.* **1** (1993) 1-23
23. Wellman, M. P.: Market-Oriented Programming: Some Early Lessons. In: Clearwater, S. H. (ed.): *Market-Based Control*. World Scientific Publishing, Singapore (1996)