

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
27 January 2011 (27.01.2011)

(10) International Publication Number
WO 2011/009652 A2

(51) International Patent Classification:
G06F 9/50 (2006.01)

00144 Roma (IT). **GIANFAGNA, Leonida** [IT/IT]; IBM Italy, Via Sciangai 53, I-RM 00144 Roma (IT).

(21) International Application Number:
PCT/EP2010/056727

(74) Agent: **BELL, Mark**; Cie IBM France, Département de la Propriété Intellectuelle, Le Plan du Bois, F-06610 La Gaude (FR).

(22) International Filing Date:
17 May 2010 (17.05.2010)

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PE, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
09165946.6 21 July 2009 (21.07.2009) EP

(71) Applicant (for all designated States except US): **INTERNATIONAL BUSINESS MACHINES CORPORATION** [US/US]; New Orchard Road, Armonk, New York 10504 (US).

(71) Applicant (for MG only): **COMPAGNIE IBM FRANCE** [FR/FR]; 17 Avenue de l'Europe, F-92275 Bois-Colombes Cedex (FR).

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

(72) Inventors; and

(75) Inventors/Applicants (for US only): **DELLA CORTE, Gianluca** [IT/IT]; IBM Italy, Via Sciangai 53, I-RM 00144 Roma (IT). **BORGHETTI, Stefano** [IT/IT]; IBM Italy, Via Sciangai 53, I-RM 00144 Roma (IT). **SGRO', Antonio** [IT/IT]; IBM Italy, Via Sciangai 53, I-RM

[Continued on next page]

(54) Title: A METHOD AND SYSTEM FOR JOB SCHEDULING IN DISTRIBUTED DATA PROCESSING SYSTEM WITH IDENTIFICATION OF OPTIMAL NETWORK TOPOLOGY

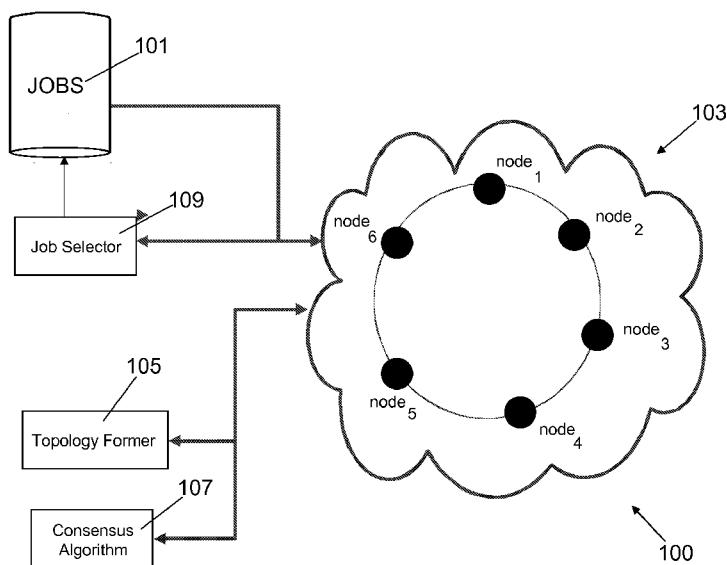


Fig. 1

(57) Abstract: The method of the present invention provides an automatic and optimised selection of the network topology for distributing scheduling of jobs on the computers of the modified network topology. The automatic and optimised selection of the network topology starts from the current topology and a desired number of additional connections. In this way the method of the present invention provides a higher convergence speed for the modified consensus algorithm in comparison e.g. to a simple ring network. The method exploits the so called small-world networks. Small-world networks are more robust to perturbations than other network architectures. The preferred embodiment provides a workload scheduling system which is highly scalable to accommodate increasing workloads within a heterogeneous distributed computing environment. A modified average consensus algorithm is used to distribute network traffic and jobs amongst a plurality of computers.

WO 2011/009652 A2

Declarations under Rule 4.17:

— *of inventorship (Rule 4.17(iv))*

Published:

— *without international search report and to be republished upon receipt of that report (Rule 48.2(g))*

5

10

**A METHOD AND SYSTEM FOR JOB SCHEDULING IN DISTRIBUTED
DATA PROCESSING SYSTEM WITH IDENTIFICATION OF OPTIMAL
NETWORK TOPOLOGY**

15

TECHNICAL FIELD

The present invention relates to the field of computer network, more particularly of job scheduling systems in a distributed computing environment.

20

BACKGROUND OF INVENTION

Job scheduling and workload balancing among a plurality of resources connected with a network is an increasingly important component of an IT environment. Many grid computing environments are driven by the scheduling of work across a distributed set of resources (e.g. computation, storage, communication capacity, software licenses, special equipment etc.). In essence, scheduling is an optimization problem, which is fairly straightforward when only one resource type is involved. However, whilst further performance improvements can be achieved by including more resource variables in the

30

scheduling process, the resulting multivariate optimization becomes a difficult mathematics problem.

State of the art job scheduling systems normally employ a master/agent architecture, wherein jobs are set up, scheduled and administered from a central server (known as a "master" server). The actual work is done by agents installed on the other servers. In use, the master maintains and interprets information relating to the jobs, available servers etc., so as to decide where to assign jobs. The agents, in turn, await commands from the master, execute the commands, and return an exit code to the master. Whilst, the master/agent architecture allows tight control over jobs, the need for the master and agents to remain synchronised (and corresponding dependency on the availability of the network and the master) is a serious limitation of the architecture. In a related manner, the highly-centralized nature of network traffic between the master and agents can degrade the overall performance of the architecture. Another problem is the limited scalability of the master/agent architecture. In particular, a master can support only a limited number of agents and creating a new master or instance creates a new and separate administration, so that the more instances created, the more complex administration activities become.

European Patent Application no. 08154507.1 filed on 15 April 2008 by the same Applicant discloses a workload scheduling system which is highly scalable to accommodate increasing workloads within a heterogeneous distributed computing environment. More particularly, the preferred embodiment employs a modified average consensus algorithm to evenly distribute network traffic and jobs amongst a plurality of computers. State information from each

computer is propagated to the rest of the computers by the modified average consensus algorithm, thereby enabling the preferred embodiment to dispense with the need for a master server, by allowing the individual
5 computers to themselves select jobs which optimally match a desired usage of their own resources to the resources required by the jobs. A drawback of the above method is that the user should establish a virtual network comprising a logical topology of the computers. In other
10 words it is the user having the responsibility of selecting the right topology and this can bring to a wrong selection which jeopardize the efficiency of the network.

It would be desirable for the user being guided in this
15 selection process or even better being able to count on a reliable method which determines the best solution according to predetermined parameters.

It is an object of the present invention to provide a
20 technique which alleviates the above drawback of the prior art.

SUMMARY OF THE INVENTION

In a preferred embodiment, the present invention provides
25 a method, in a network of a plurality of N computers being connected together, for executing jobs based on computation of an improved network topology, the improved network topology including a number C of additional connections with respect to the current topology, the
30 method including the steps of representing the current topology of the network of the plurality of N computers by means of $N \times N$ Laplacian matrix M wherein the values $l_{i,j}$ are defined as follows

$$l_{i,j} := \begin{cases} \text{deg}(v_i) & \text{if } i = j \\ -1 & \text{if } i \neq j \text{ and } v_i \text{ adjacent } v_j \\ 0 & \text{otherwise} \end{cases}$$

wherein $\text{deg}(v_i)$ is the number of nodes that are connected to the node i and wherein $1 \leq i \leq N$ and $1 \leq j \leq N$; then the
 5 $N \times N$ Laplacian matrix M_1 is calculated, having

$\text{Tr}(M_1) = \text{Tr}(M) + C \cdot 2$ and having the greatest possible second-smallest eigenvalue; the topology is thus modified into an improved topology by adding C additional connections according to the calculated $N \times N$ Laplacian

10 matrix M_1 ; then, starting from the modified network topology for distributing scheduling of jobs on the computers of the modified network topology.

The method of the present invention can help to
 15 solve the problem of the prior art by providing an automatic and optimised selection of the network topology starting from the current topology and a desired number of additional connections. In this way the method of the present invention provides a higher convergence speed for
 20 the modified consensus algorithm in comparison e.g. to simple ring network. The method exploits the so called small-world networks. Small-world networks are more robust to perturbations than other network architectures, in fact in a random network, in which all nodes have
 25 roughly the same number of connections, deleting a random node is likely to increase the mean-shortest path length slightly but significantly for almost any node deleted. In this sense, random networks are vulnerable to random perturbations, whereas small-world networks are robust.

In a further embodiment of the present invention it is provided a system comprising components adapted to implement the method above.

5 In another embodiment a computer program is provided which realizes the method above when run on a computer.

BRIEF DESCRIPTION OF THE DRAWINGS

10

Embodiments of the invention will now be described, by way of example only, by reference to the accompanying drawings, in which:

15

Figure 1 is a block diagram of a software architecture of a workload scheduling system of the preferred embodiment with a representation of a ring topology;

20

Figure 2 shows the topology of Figure 1 modified according to the method of the present invention;

Figure 3a and 3b show an example of a computer network to which the method of the present invention is applied;

25

Figure 4 is a diagram of a general computer system adapted to support the method of workload scheduling of the preferred embodiment; and

Figure 5 shows a flowchart representing the steps to perform a method according to a preferred embodiment of the present invention.

30

DESCRIPTION OF A PREFERRED EMBODIMENT

As shown in Figure 1, the preferred embodiment 100 comprises a repository 101 of jobs to be performed in a distributed heterogeneous network 103.

For simplicity, assume that the distributed network 103 comprises n nodes wherein each node ($node_i$) possesses a number of resources. The preferred embodiment employs a topology-forming algorithm 105 which will be discussed in deeper details; a modified average consensus algorithm 107 is used to enable nodes in the distributed network 103 to advise other nodes in the network of their current status and availability to execute new jobs. The preferred embodiment further comprises a job-selection module 109, which enables the nodes in the distributed network 103 to select an optimal job (job_k) (to execute next) from the job repository 101, in accordance with a user-defined desired usage of the resources of each node ($node_i$) and the resource requirements of the job (job_k).

The topology forming algorithm 105 establishes a virtual network comprising a logical topology of the nodes in the distributed heterogeneous network 103. Within, the virtual network, the logical topology establishes which nodes can communicate with each other. In particular, the logical topology is defined so that each node is directly connected to (and can communicate with) j neighbouring nodes (wherein $j < n-1$). Thus, $node_i$ is provided with a neighbourhood N_i , comprising j nodes ($node_p$, $p=1$ to j). The number of nodes to which a node is connected can be defined by the user, wherein the fault-tolerance (and convergence rate) of the preferred embodiment is improved by increasing the number of such nodes. In a preferred embodiment of the present invention, starting from a basic topology e.g. the network 103 of Fig 1 the user can indicate a number of

additional connections according to the monitored performances of the network and the traffic.

Alternatively the number of optimal additional

connections could be determined by means of monitoring

5 tools which are able to monitor the network performances in order to evaluate the efficiency and speed of the

network; another possibility is to rely on statistics of previously measured performances. The only condition

imposed on the basic topology is that the graph formed by

10 the virtual network must be connected, in other words,

starting from each node it must be possible to reach any other node through an arbitrary number of steps.

In a ring topology as the one showed in Fig 1 each

node has only two connections and two neighbours. So,

15 for example, referring to Figure 1, node₁ is directly

connected to node₂ and node₆. Similarly, node₂ is directly

connected to node₁ and node₃; and node₆ is directly

connected to node₁ and node₅. Any topology of n nodes can

be represented with a n x n matrix. The example of Fig. 1

20 can be represented with the following 6x6 matrix M:

$$\begin{matrix}
 & 2 & -1 & 0 & 0 & 0 & -1 \\
 & -1 & 2 & -1 & 0 & 0 & 0 \\
 25 & 0 & -1 & 2 & -1 & 0 & 0 \\
 & 0 & 0 & -1 & 2 & -1 & 0 \\
 & 0 & 0 & 0 & -1 & 2 & -1 \\
 & -1 & 0 & 0 & 0 & -1 & 2
 \end{matrix}$$

30 The Eigenvalues of such matrix are:

0, 1, 1, 3, 3 and 4

As explained by Yoonsoon Kim and Mehran Mesbahi (see Yoonsoon Kim, Mehran Mesbahi - "On maximizing the Second

Smallest Eigenvalue of a State dependent Graph Laplacian"

- IEEE TRANSACTIONS ON AUTOMATIC CONTROL, VOL.51, No 1,
JANUARY 2006) given a graph G with n nodes (without
loops or multiple edges), its Laplacian matrix L is

5 defined as:

$$l_{i,j} := \begin{cases} \text{deg}(v_i) & \text{if } i = j \\ -1 & \text{if } i \neq j \text{ and } v_i \text{ adjacent } v_j \\ 0 & \text{otherwise} \end{cases}$$

where $\text{deg}(v_i)$ is the number of nodes that are
connected to the node i . Starting from a current topology
10 (e.g. the Laplacian matrix M represented above) and
assuming to add C additional connections among the
existing nodes, the problem to be solved is to find the
Laplacian matrix M_1 having the trace (Tr) equal to the
trace of matrix $M + C \cdot 2$ with the greatest possible
15 second-smallest eigenvalue;

In other words:

determining among all possible Laplacian matrix
wherein $\text{Tr}(M_1) = \text{Tr}(M) + C \cdot 2$ that matrix M_1 having the
greatest possible second-smallest eigenvalue.

20

On the cited article by Yoonsoon Kim and Mehran
Mesbahi it is demonstrated that using the proposed
protocol the convergence speed is proportional to the
second smallest eigenvalue of the Laplacian related to
25 the network graph.

So in order to maximize the convergence speed we can use
the approach proposed by Yoonsoon Kim and Mehran Mesbahi,
in order to obtain the best logical topology to solve the
consensus problem. Using this method it is possible to
30 obtain a logical topology that is the best small-world
network.

In mathematics and physics, a small-world network is a type of mathematical graph in which most nodes are not neighbours of one another, but most nodes can be reached from every other by a small number of hops or steps.

5 In this way it is possible to decrease the mean-shortest path length, so the states propagation is faster.

The algorithm by Yoonsoon Kim and Mehran Mesbahi (Kim-Mesbahi algorithm) is initiated at time $t=0$ with an
 10 initial graph (configuration) G_0 and then for $t=0, 1, 2, \dots$ we proceed to iteratively find the graphs that maximises $\lambda_2(LG(t+1))$ where $LG(t+1)$ is the Laplacian of graph G at time $t+1$. This greedy procedure is then iterated upon until the value of $\lambda_2(LG(t))$
 15 can not be improved further. We note that the proposed greedy algorithm converges, as the sequence generated it is non-decreasing and bounded.

With the method described by Yoonsoon Kim and Mehran Mesbahi it is possible, starting from a current topology
 20 and indicating how many additional connections are requested, to identify the best topology in order to maximise the convergence speed, i.e. the second smallest eigenvalue of the matrix representing the topology. Applying this algorithm to the ring topology represented
 25 with the matrix above we would obtain the following new matrix:

	3	-1	0	-1	0	-1
	-1	3	-1	0	-1	0
	0	-1	3	-1	0	-1
30	-1	0	-1	3	-1	0
	0	-1	0	-1	3	-1
	-1	0	-1	0	-1	3

which corresponds to the topology represented in Figure 2.

The eigenvalues of the new matrix are the following:
5 0, 3, 3, 3, 3 and 6 which means that the convergence speed is proportional to 3 (the second smallest eigenvalue).

It will of course be realised that the topology
10 shown in Figure 1 and Figure 2 is provided for example purposes only and should in no way be construed as limiting the preferred embodiment to a ring topology. In particular, the skilled person will understand that the preferred embodiment is operable with any topology and
15 number of nodes in the distributed network 103 (subject to the above-mentioned connected constraint). Figure 3a shows another (more complex) example. In this basic topology 26 nodes are connected in a ring configuration. Such a configuration is represented in a matrix of 26x26
20 (not represented here) having the second smallest eigenvalue equal to 0.0581. If we want to add 6 more connections to such topology and we apply the Kim-Mesbahi algorithm we would obtain the topology shown in Figure 3b which has an eigenvalue equal to 0.2413: this corresponds
25 to a consensus convergence speed 415% higher than previous topology configuration and this is the highest value that is possible to obtain adding 6 connections.

Each node ($node_i$) in the virtual network comprises a
30 used resources state variable indicating the extent to which the nodes' resources are occupied by the jobs currently running thereon. This information can be acquired from real-time resource consumption metrics gathered by monitoring software. In particular, defining

the \underline{o}_i as the consumption metric vector associated with a given resource vector \underline{res}_i , the extent to which a given node's resources are occupied (by currently running jobs) can be given by a scalar variable $state_i$, which is a weighted sum of the consumption metrics of the resource variables of a node ($node_i$) (i.e. $state_i = \sum_{l=1}^r a_{li} o_{li}$ [or in vector notation, $\underline{state} = O^T \text{diag}(a)$. (wherein $\underline{state} \in \mathbb{R}^{nx1}$, $O \in \mathbb{R}^{rxn}$ and $\underline{a} \in \mathbb{R}^{rxr}$). The modified average consensus algorithm 16 in the preferred embodiment enables the nodes to propagate this information throughout the entire virtual network even to nodes to which the originating node is not directly connected.

A similar notation may be used to describe the requirements of a particular job (in the job repository). In particular, the net requirements (e_q) of a job (job_q) may be defined as $e_q = \sum_{l=1}^r \beta_{ql} req_{ql}$.

As explained in European Patent Application no. 08154507.1 the job scheduler is based on a modified version of the so called Consensus algorithm.

The average consensus model of a graph provides a distributed method of calculating graph evolution with an input u_i to a node i and its neighbours (N_i). Thus, if the evolution of the $state_i$ of node i can be denoted by

$$state'_i = f_i(state_i) + u_i(state_i, state_p), p = 1 \text{ to } j, j \in N_i \Rightarrow state'_i = u_i(state_i, state_p),$$

it can be demonstrated that $state'_i = - \sum_{p=1}^{|N_i|} (state_i - state_p)$

(continuous solution) or

$$state_i(k+1) = \frac{1}{|N_i|+1} \left(state_i(k) + \sum_{p=1}^{|N_i|} state_p(k) \right) \quad (\text{discrete solution})$$

asymptotically solves the consensus problem into a connected graph. Accordingly, the dynamic system

converges to the mean of the initial states,

$$\lim_{k \rightarrow \infty} state_i = \frac{1}{j} \sum_{p=1}^j state_p(0) \quad (\text{continuous solution}) \quad \text{or}$$

$$\lim_{k \rightarrow \infty} state_i = \frac{1}{j} \sum_{p=1}^j state_p(0) \quad (\text{discrete solution}),$$

wherein the proof for these limits is derived from the related Nyquist

5 diagram therefore.

The preferred embodiment modifies the above-mentioned traditional average consensus algorithm by introducing a virtual node V (not shown), which is directly connected to all of the other nodes within the virtual network. Thus, the virtual node V is included within the neighbourhood N_i of a node i . Using this approach, each node (node i) calculates its next state ($state_i(k+1)$) from:

- 15 - its current state ($state_i(k)$); and
- the current states of the other nodes (including the virtual node V) in its neighbourhood N_i using the

following formula
$$state_i(k+1) = \frac{1}{|N_i|+1} \left(state_i(k) + \sum_{p=1 \text{ to } j, p \in N_i} state_p(k) \right).$$

The virtual node V has a user-configurable, fixed state which represents the desired workload of all the nodes in the virtual network. The inclusion of the virtual node V into the neighbourhood of each node in the virtual network causes the average consensus algorithm with all the other nodes in the virtual network to balance against and converge to the fixed state of the virtual node V (wherein the convergence speed is related to the Laplacian of the network graph). Thus, by making the state of the virtual node V configurable by the user, the preferred embodiment effectively provides a mechanism for tuning a workload schedule to meet a desired usage of the resources of the nodes (i.e. operating point) in the

30

virtual network (i.e. to effectively alter the operating point of the virtual network).

The job-selection module 109 enables a node ($node_i$)
5 to select a job (from the job repository 10) to execute next, in accordance with the node's current state and its calculated next state. In particular, if $state_i(k+1) < state_i(k)$, then no new job is to be undertaken by the node at the next iteration. However, *if $state_i(k+1) \geq state_i(k)$* , a
10 difference variable Δ is defined as $\Delta = state_i(k+1) - state_i(k)$ (i.e. the difference between the calculated next state of $node_i$ and the current state of the node). The next job (job_{k+1}) selected (from the job repository) is the job (job_t) whose net requirements variable (e_t) has minimal
15 difference from the difference variable Δ (i.e. $e_t | \min_t (\Delta - e_t)$).

With reference to Figure 4 a generic computer of the system (e.g. computer, Internet server, router, remote
20 servers) is denoted with 450. The computer 450 is formed by several units that are connected in parallel to a system bus 453. In detail, one or more microprocessors 456 control operation of the computer 450; a RAM 459 is
25 directly used as a working memory by the microprocessors 456, and a ROM 462 stores basic code for a bootstrap of the computer 450. Peripheral units are clustered around a local bus 465 (by means of respective interfaces). Particularly, a mass memory consists of a hard-disk 468 and a drive 471 for reading CD-ROMs 474. Moreover, the computer 450
30 includes input devices 477 (for example, a keyboard and a mouse), and output devices 480 (for example, a monitor and a printer). A Network Interface Card 483 is used to connect the computer 450 to the network. A bridge unit 486

interfaces the system bus 453 with the local bus 465. Each microprocessor 456 and the bridge unit 486 can operate as master agents requesting an access to the system bus 453 for transmitting information. An arbiter
5 489 manages the granting of the access with mutual exclusion to the system bus 453. Similar considerations apply if the system has a different topology, or it is based on other networks. Alternatively, the computers have a different structure, include equivalent units, or
10 consist of other data processing entities (such as PDAs, mobile phones, and the like).

Figure 5 schematically shows the method steps according to a preferred embodiment of the present invention (500).
15 The process starts at step 501 and goes to step 503 where the current topology of the network is determined and represented with a matrix M as explained above. According to a preferred embodiment of the present invention, the dimension of the matrix M is $n \times n$ where n is the number
20 of the nodes in the network. At step 503 the user (or administrator) of the network is allowed to indicate how many additional connections are desired. The determination of the number of required additional connections can take into account several parameters
25 (e.g. network traffic, system performance, geographic configuration of the network) and can be determined in several different ways; those skilled in the art will appreciate that this determination can be done with existing tools or even indicated manually by the
30 administrator. The number of optimal additional connections could also be determined by means of monitoring tools which are able to monitor the network performances in order to evaluate the efficiency and speed of the network; another possibility is to rely on

statistics of previously measured performances. The matrix M and the number of additional connections required are then input to the Kim-Mesbahi which is able to determine the optimal topology of the network, by
5 finding the matrix M_1 having the maximum second-smallest eigenvalue of all possible solutions. In a preferred embodiment of the present invention the process continues as explained above with the job scheduling process (step 509) which is based on the modified consensus algorithm.
10 Optionally (step 511) the performances of the new network are monitored to determine whether they are satisfactory or not. If they are determined to be satisfactory the normal job scheduling activities can continue, otherwise a new re-configuration of the network is considered and
15 the control goes back to step 505. Another possible alternative could be to increment by 1 the number of additional connections until the performances are determined to be satisfactory at step 511.

20 Alterations and modifications may be made to the above without departing from the scope of the invention. Naturally, in order to satisfy local and specific requirements, a person skilled in the art may apply to the solution described above many modifications and
25 alterations. Particularly, although the present invention has been described with a certain degree of particularity with reference to preferred embodiment(s) thereof, it should be understood that various omissions, substitutions and changes in the form and details as well
30 as other embodiments are possible; moreover, it is expressly intended that specific elements and/or method steps described in connection with any disclosed embodiment of the invention may be incorporated in any other embodiment as a general matter of design choice.

For example, similar considerations apply if the computers have different structure or include equivalent units; in any case, it is possible to replace the computers with any code execution entity (such as a PDA,
5 a mobile phone, and the like).

Similar considerations apply if the program (which may be used to implement each embodiment of the invention) is structured in a different way, or if additional modules or functions are provided; likewise,
10 the memory structures may be of other types, or may be replaced with equivalent entities (not necessarily consisting of physical storage media). Moreover, the proposed solution lends itself to be implemented with an equivalent method (having similar or additional steps,
15 even in a different order). In any case, the program may take any form suitable to be used by or in connection with any data processing system, such as external or resident software, firmware, or microcode (either in object code or in source code). Moreover, the program may
20 be provided on any computer-usable medium; the medium can be any element suitable to contain, store, communicate, propagate, or transfer the program. Examples of such medium are fixed disks (where the program can be pre-loaded), removable disks, tapes, cards, wires, fibres,
25 wireless connections, networks, broadcast waves, and the like; for example, the medium may be of the electronic, magnetic, optical, electromagnetic, infrared, or semiconductor type.

In any case, the solution according to the present
30 invention lends itself to be carried out with a hardware structure (for example, integrated in a chip of semiconductor material), or with a combination of software and hardware.

CLAIMS

1. A method, in a network of a plurality of N computers
 5 being connected together, for executing jobs based on
 computation of an improved network topology, the improved
 network topology including a number C of additional
 connections with respect to the current topology, the
 method including the steps of:

10 - representing the current topology of the network
 of the plurality of N computers by means of NxN Laplacian
 matrix M wherein the values $l_{i,j}$ are defined as follows

$$l_{i,j} := \begin{cases} \text{deg}(v_i) & \text{if } i = j \\ -1 & \text{if } i \neq j \text{ and } v_i \text{ adjacent } v_j \\ 0 & \text{otherwise} \end{cases}$$

15 wherein $\text{deg}(v_i)$ is the number of nodes that are
 connected to the node i and wherein $1 \leq i \leq N$ and $1 \leq j \leq N$;

- calculating the NxN Laplacian matrix M_1 having
 $\text{Tr}(M_1) = \text{Tr}(M) + C \cdot 2$ and having the greatest possible
 second-smallest eigenvalue;

20 - modifying the network topology into an improved
 topology by adding C additional connections according to
 the calculated NxN Laplacian matrix M_1 ;

- starting from the modified network topology for
 distributing scheduling of jobs on the computers of the
 25 modified network topology.

2. The method of claim 1 wherein the improved network
 topology is a small-world network.

3. The method of any preceding claims wherein the step of calculating the NxN Laplacian matrix M_1 having $\text{Tr}(M_1) = \text{Tr}(M) + C \cdot 2$ and having the greatest possible second-smallest eigenvalue includes applying the Kim-Mesbahi

5 algorithm.

4. The method of any preceding claim further comprising the step of:

10 - prompting a user for input a desired value of the number of C additional connections.

5. The method of any preceding claim further comprising the steps of:

15 - monitoring the network performances to measure a value indicative of the efficiency of the network;
- determining the number C of additional connections according to the measured efficiency.

6. The method of any one of claims 1 to 5 wherein the step of distributing scheduling of jobs on the computers of the modified network topology comprises:

20 - establishing a desired at least one operational resource value for the network inside the network topology;
25 - determining the current usage of the resources of at least some of the computers in the network, by one or more jobs being executed thereon;
- calculating a predicted state value for each computer in the network from the current usage of the computers resources and the desired operating point; and
30 - selecting another job to be executed next by one of the computers in the network in the event the

computer's predicted state value substantially exceeds the current usage of the computer's resources.

7. The method as claimed in claim 6 wherein the step of
5 determining the current usage of the resources of at least some of the computers in the network comprises the step of acquiring information relating to the current usage of the resources from resource consumption metrics gathered by monitoring software.

10

8. Method as claimed in claim 6 or 7 wherein in the event a computer's predicted state value substantially exceeds the current usage of the computer's resources, a job is selected for execution thereby, whose resource
15 requirements are closest to the difference between the computer's predicted state value and the current usage of the computer's resources.

9. System for scheduling a workload for a plurality of
20 computers wherein the system comprises a one or more components adapted to perform the method of any claim 1 to 8.

10. Computer program comprising instructions for carrying
25 out the method of any claim 1 to 8 when said computer program is executed on a computer system.

11. A service deployed in a data processing system for implementing the method of any claim 1 to 8.

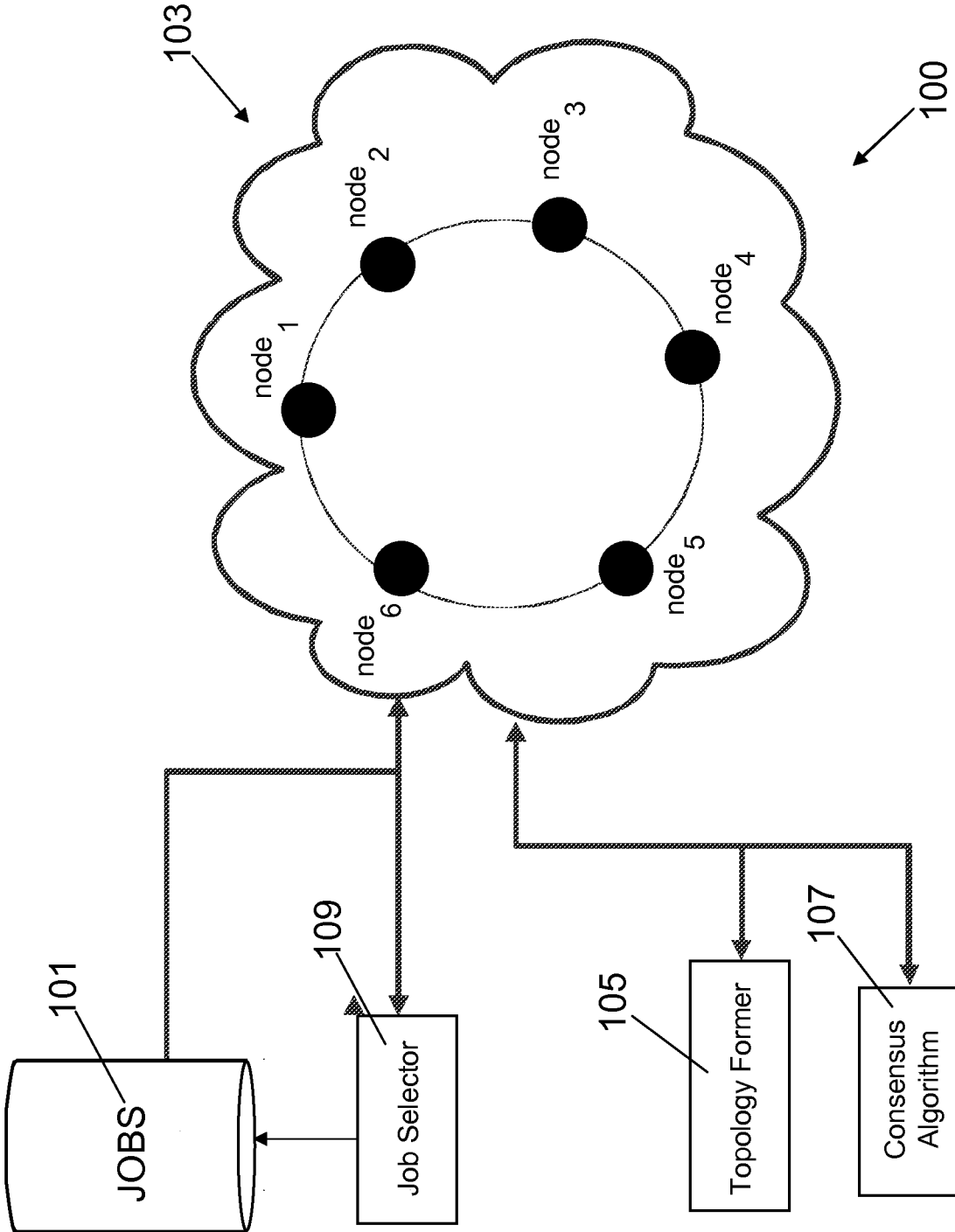


Fig. 1

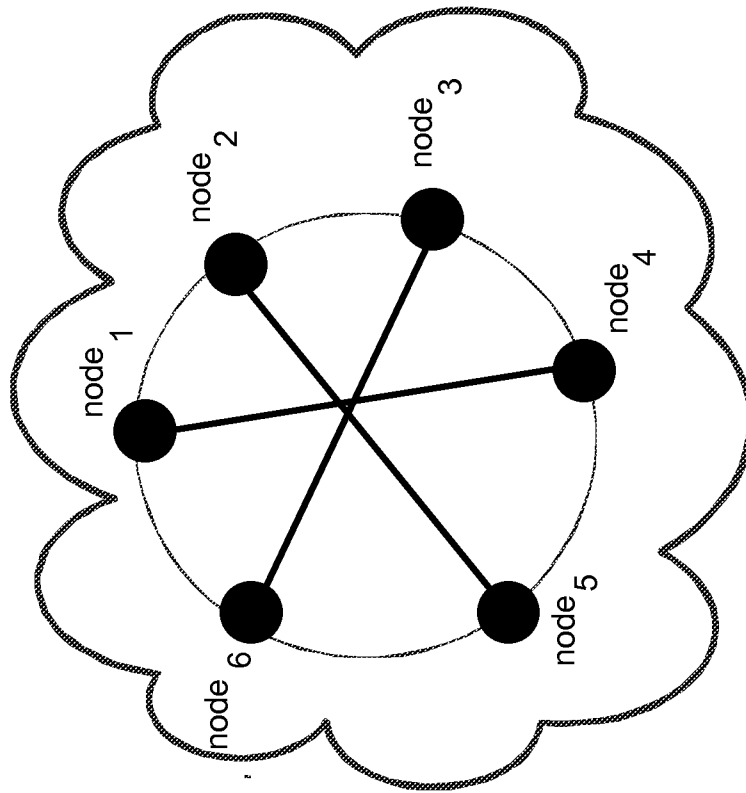


Fig. 2

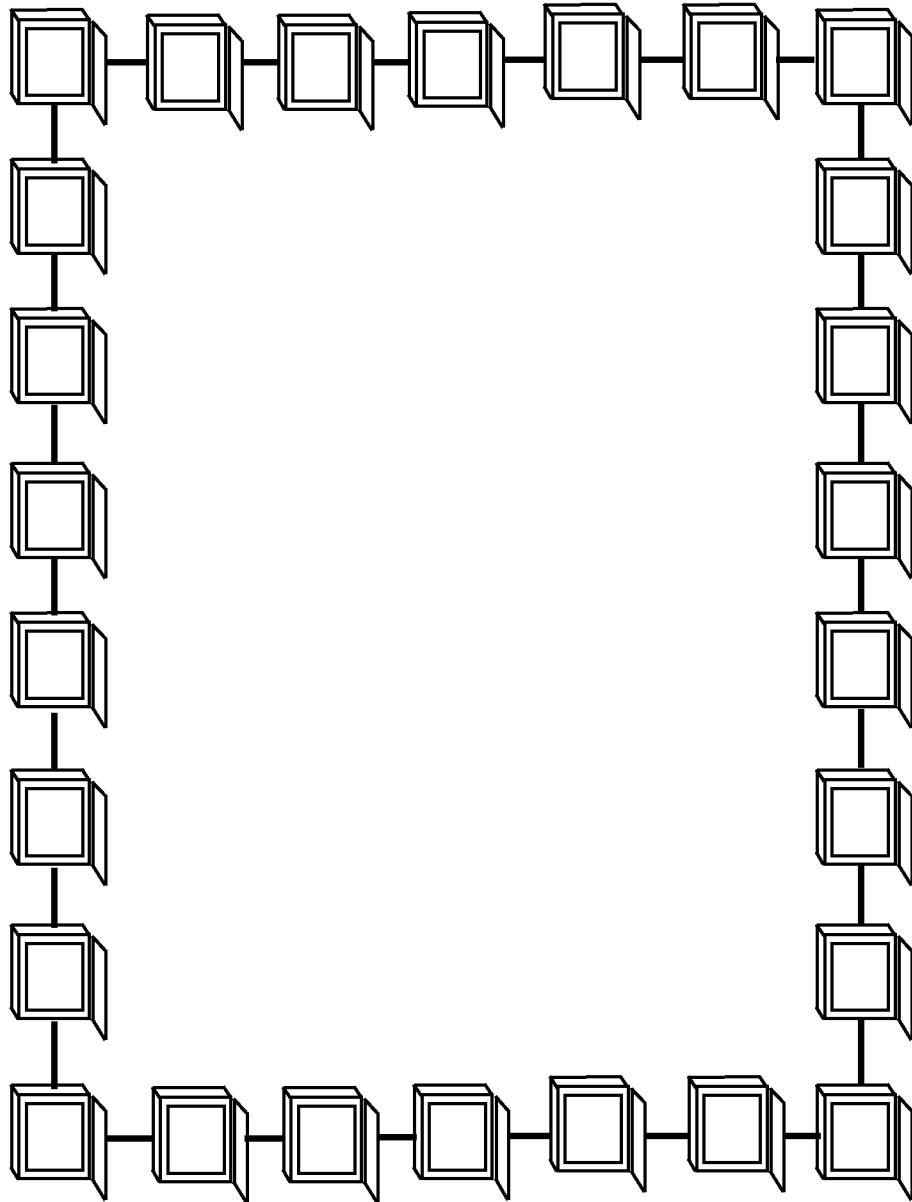


Fig. 3A

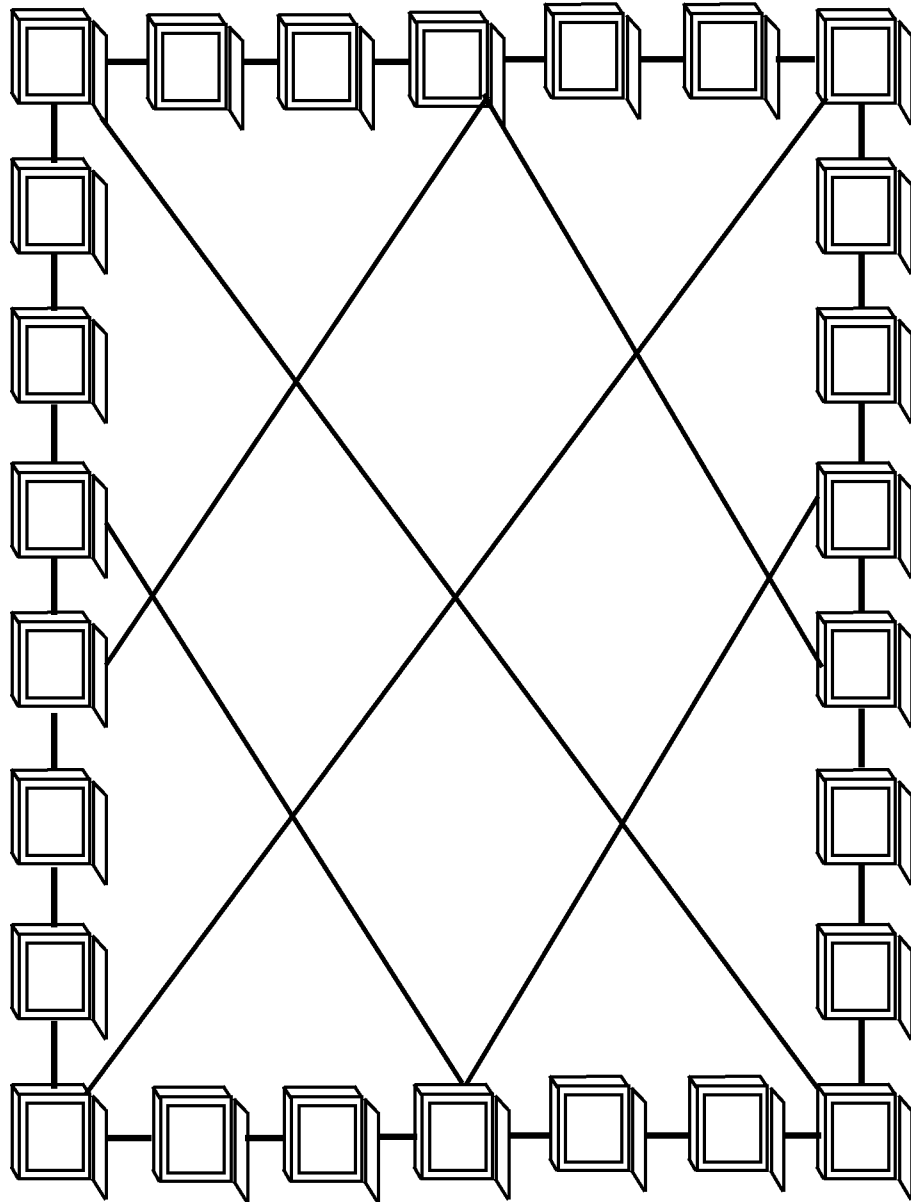


Fig. 3B

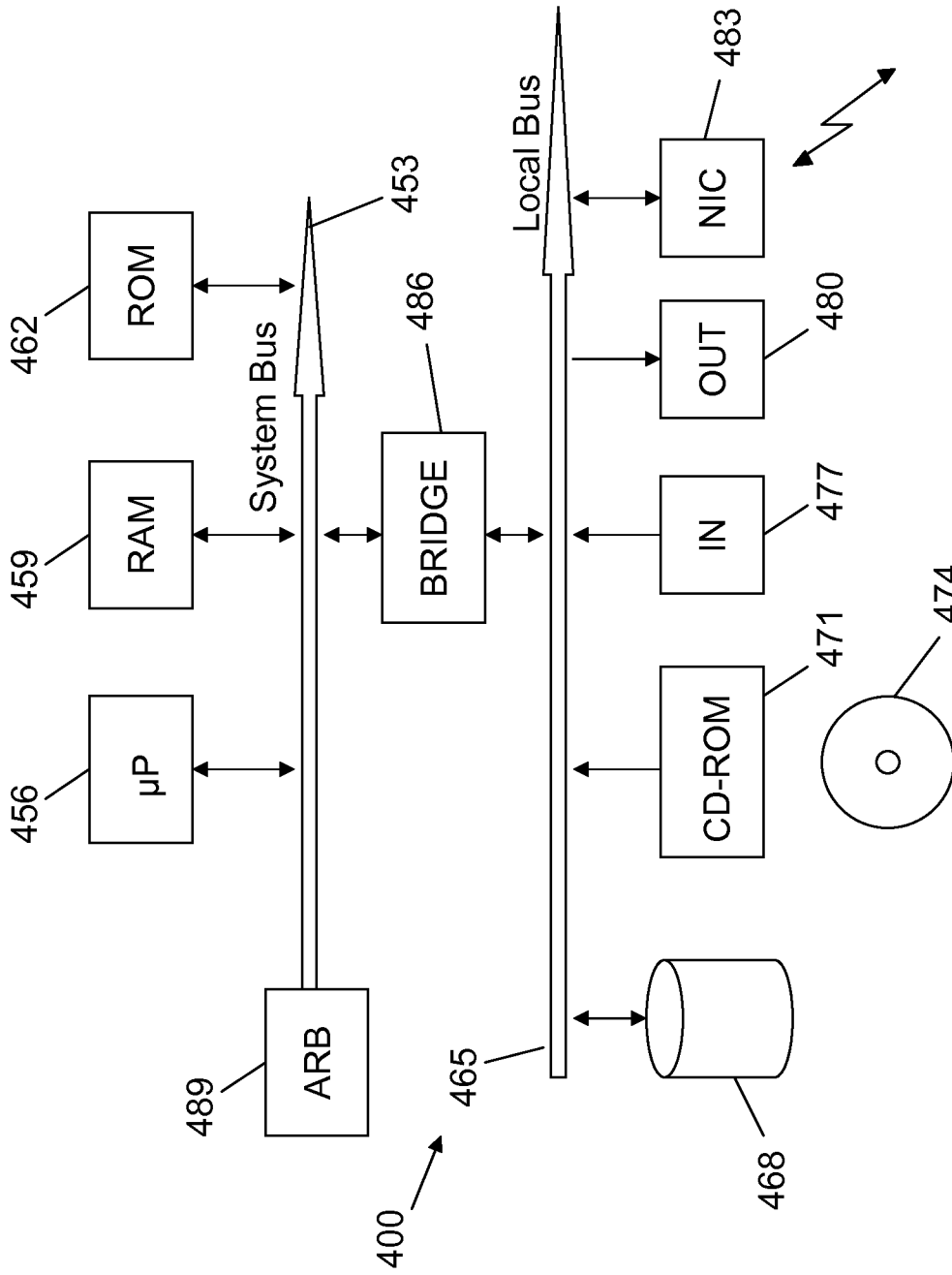


Fig. 4

