

An Efficient Load Balancing Algorithm for P2P Systems

Khaled Ragab

Computer Science Dept., College of Computer Science and Information Technology,
Hofuf, Saudi Arabia
kabdultawab@kfu.edu.sa

Abstract—Peer-to-Peer (P2P) networks are considered to be the most important development for content distribution and sharing infrastructures. The quality of services provided by the majority of the P2P systems is questioned. Load balancing among peers is critical and a key challenge. This paper proposes a load balancing algorithm based on stochastic analysis. It addresses the out-of-date problem as a result of node's state changes during loads movement among nodes.

Index Terms— Peer-to-Peer, Overlay Networks, load balancing.

I. INTRODUCTION

In recent years the development of distributed systems especially the internet has been influenced heavily by the **Peer-to-Peer (P2P)** paradigm. P2P is a class of applications that takes advantage of resources e.g., storage, cycles, content, human presence, available at the edges of the Internet [1]. P2P systems offer an alternative to such traditional client-server systems for several application domains. They have emerged as an interesting solution for content distribution, sharing and locating resources over the Internet. In P2P systems, every node (peer) of the system acts as both client and server (**servant**) and provides part of the overall resources/information available from the system. Each node often has different resource capabilities (e.g. processor, storage, and bandwidth) [2]. Thus, it is required that each node has a load proportional to its resources capabilities. On account of the dynamism natures of the P2P systems, it is difficult to ensure that the load is uniformly distributed across the system. In particular, this paper considers a P2P system of M nodes in which nodes join/leave, and data entity inserted/deleted continuously. Similarly to [3], [4], [5] this paper assumes node and data entity have been assigned identifiers that chosen randomly. Thus, there is a $\Theta(\log M)$ imbalance factor in the number of data entities stored at a node. Additionally, the imbalance factor becomes more worse if the P2P applications associate semantics with data entity IDs since IDs will

not be uniformly distributed. Consequently, it is important to design mechanisms that balance the system load. There are two distinct strategies to distribute the system workload [6]. First, load balancing algorithms that strive to equalize the workload among nodes. Second, load sharing algorithms which simply attempt to assure that no node is idle while jobs at other nodes are waiting for service.

Load balancing techniques in P2P systems should be scalable and cope with its large size. They should place or re-place shared data entities optimally among nodes while maintaining an efficient overlay routing tables to redirect queries to the right node.

The communication delays among peers significantly alter the expected performance of the load balancing schemes. Due to such delay, the information that a particular peer has about other peers at any time is dated and may not accurately represent the current state of the other peers. For the same reason, a load sent to a recipient peer arrives at a delayed instant. In the mean time, however, the load state of the recipient peer may have considerably changed from what was known to the transmitting peer at the time of load transfer. This paper proposes a stochastic dynamic load balancing algorithm that tackles the out-of-date problem.

The remainder of this paper is organized as follows. Section two introduces a survey for the load balancing algorithms. Section three exposes the proposed stochastic load balancing model and algorithm. Evaluation of the proposed has been discussed in section four. Section five draws a conclusion of this paper.

II. LOAD BALANCING SURVEY

Load balancing is the problem of mapping and remapping workload in the distributed system.

A. Load Balancing Design

Load balancing design determines how nodes communicate and migrate loads for the purpose of load balancing. It moves workload from heavily loaded nodes (senders) to lightly loaded nodes (receivers) to improve the system overall performance [15]. Load balancing design includes four components that can be classified as follows [16] [17].

- *Transfer policy*: It decides whether a node is in a suitable state to participate in a load transfer; either

Manuscript received February 15, 2011; revised May 15, 2011; accepted July 15, 2011.

The author is on leave holiday from Ain Shams University, Cairo, Egypt.

receiver or sender.

- *Location policy*: Once the transfer policy decides that a node is a receiver or sender. The location policy takes the responsibility to find a suitable sender or receiver.
- *Selection policy*: Once the transfer policy decides that a node is a sender, the selection policy specifies which load should be transferred. It should take into account several factors such as load transfer cost, and life time of the process that should be larger than load transfer time.
- *Information policy*: It decides when and how to collect system state information.

Load balancing designs are categorized into static and dynamic. With a static load balancing scheme, loads are scattered from sender to receiver through deterministic splits. Static schemes are simple to implement and easy to achieve with slight overhead [13]. They perform perfectly in homogenous systems, where all nodes are almost the same, and all loads are same as well. On the other hand, the dynamic load balancing schemes make decisions based on the current status information [15]. Accordingly, the transfer policy at certain node decides to be a sender or receiver, the selection policy selects the load to be transferred. The dynamic load balancing schemes perform efficiently when its nodes have heterogeneous loads, and resources. The typical architectures of dynamic load balancing schemes can be classified into centralized, distributed, and topological. In a centralized scheme, a central server “*coordinator*” receives load report from the other nodes, while overloaded nodes request the coordinator to find underloaded nodes [17]. In distributed architecture, each node has a global or partial view of the system status. Consequently, the transfer policy at each node can locally decide to transfer a load either out from it (*sender-initiated*) or into it (*receiver-initiated*) [18]. Then, the location policy at each node probes a limited number of nodes to find a suitable receiver or sender. *P. Kruger* [19], proposed symmetrically-initiated adaptive location policy that uses information gathered during its previous searches in order to keep track of the recent state of each node in the system. It finds a suitable receiver when a heavily-loaded node wishes to send a load out, and finds a suitable sender when a lightly-loaded node wishes to receive a load. Finally, in a system with large number of nodes, a topological scheme should be used [20]. It partitions nodes into groups. The load balance is performed in each group first, then, a global load balance among groups will be performed. However nodes in the hierarchical architecture [21] are organized into a tree. Inner nodes gather the status information of its sub-trees. Then, load balancing is performed the leaves to the roots of the tree.

B. P2P Load Balancing

Load balancing is a critical issue for the efficient operation of the P2P systems. Recently, it attracted much

attention in the research community especially in the *distributed hashed table (DHT)* based P2P systems. Namespace balancing is struggling to balance the load across nodes by ensuring that each node is responsible for a balanced namespace. This is valid only under the assumption of uniform workload and uniform node capacity. Otherwise, there is a $\Theta(\log M)$ imbalance factor in the number of objects stored at a node. To mitigate this imbalance two categories, *node placement* and *object-placement* load balancing techniques were proposed. In the node placement technique, nodes can be placed or replaced in locations with heavy loads. For example, a node in the *Mercury* load balancing mechanisms [7] is able to detect a lightly loaded range, and move there if it is overloaded. In object placement technique, objects are placed at lightly loaded nodes either when they are inserted into the system [11] or through dynamic load balancing schemes based on the virtual servers (*VSs*) concept [8], whose explicit definition and use for load balance was proposed by *Godfrey et. al.* [9] and *Rao et. al.* [11]. In [8], a virtual server represents a peer in the *DHT*; that is the storage of data items. In addition, routing takes place at the virtual server level rather than at the physical node level. Each physical node has a pool of *VSs* as seen in figure 1. Load balancing could be achieved by migrating *VSs* from heavily loaded physical node to lightly loaded physical node. One main advantage of using *VSs* for balancing the load is that approach does not require any changes to the underlying *DHT*. In fact, the transfer of a virtual server can be implemented simply as a peer leaving and peer joining the system. In [11], *Rao et. al.* proposed three simple and static load balancing schemes: “*one-to-one*”, “*one-to-many*” and “*many-to-many*”. *Godfrey et. al.* combines both “*one-to-many*” and “*many-to-many*” schemes and uses them in different scenarios [9]. Clustered *VSs* scheme is presented in [12] that optimized the basic *VS* framework to reduce the overhead involved in the *VS* framework. However, *VSs* cannot be moved, and therefore, the scheme cannot respond to dynamic changes in network conditions. This paper focuses on the design and analysis of P2P load balancing algorithm based on stochastic analysis [23, 24] and based on the *VSs* concept [8].

C. Challenges: P2P Load Balancing

Load balancing techniques in P2P systems are facing challenges coming from the characteristics of these systems. First, the size of the P2P system is large that means a scalable load balancing technique is required. Second, dissimilar to the traditional systems, nodes of a P2P system are not replicas and requests cannot be executed in any node. If nodes have dated, inaccurate information about the state of other nodes, due to random communication delays between nodes, then this could result in unnecessary periodic exchange of loads among them. For example, an overloaded node removes

some of its virtual servers. However, such simple deletion will cause the problem of “load thrashing”¹, for the removed virtual servers may make other nodes overloaded. Consequently, this paper proposes a stochastic P2P load balancing algorithm that approximately determines the minimum amount of time to change the node’s state from overloaded to underloaded and vice-versa. Comparing that time with the required time to migrate virtual servers enable us to come to a careful decision. Accordingly, the proposed algorithm undoubtedly avoids the load thrashing. To the best of the author’s knowledge, there is no any load balancing algorithm for the P2P system based upon the following stochastic analysis.

III. LOAD SHARING ALGORITHM

A. Model

This paper considers a P2P system consisting of M physical nodes (peers), denoted by $P_i, 1 \leq i \leq M$. Each peer can be modeled as a queuing system, such as M/M/1, M/D/1, etc. Each physical node P_i has a capacity C_i that corresponds to the maximum amount of load that it can process per unit of time. Nodes create virtual servers (VSs), which join the P2P network. Therefore, it can own multiple noncontiguous portions of the DHT’s identifier space. Each virtual server participates in the DHT as a single entry (e.g routing table). Moreover, each virtual server stores data items whose IDs fall into its responsible region of the DHT’s identifier space. As seen in fig. 1, a node P_i might have n VSs v_1, v_2, \dots, v_n ; where $n = VSset.size$. Each v_j has load l_j ; (for $j=1, \dots, n$). The load of peer P_i in a unit of time is $L_i = l_1 + l_2 + \dots + l_n$. The utilization of a node’s P_i is L_i/C_i . From the perspective of load balancing, a virtual server represents certain amount of load (e.g. the load generated by serving the requests of the data items whose IDs fall into its responsible region) [25]. To avoid fluctuations in workload nodes should operate below their capacity. If a node finds itself receiving more load L_i than the upper target load U (i.e. $(L_i/C_i) > U$), it considers itself overloaded. A node P_i also has load L_i less than L is considered to be underloaded. An overloaded node is not able to store objects given to it route packets, or carry out computation, depending on the application.

Definition 1.

A node P_i is in one of the following state as follows

$$S_i = \begin{cases} \text{Underloaded} & \text{if } Q_i < L \\ \text{Normal} & \text{if } L \leq Q_i \leq U \\ \text{Overloaded} & \text{if } U < Q_i \end{cases}$$

Clearly the state space Q_i consists of non-negative integers sub-divided into three disjoint regions $[0, L)$, $[L, U]$, and (U, ∞) corresponding to underloaded, normal, and overloaded state respectively.

¹ Load thrashing is a condition when the load balancing algorithm is engaged in moving virtual servers back and forth between nodes.

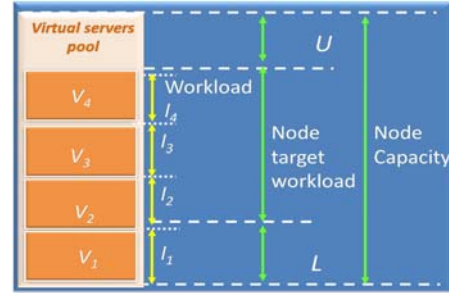


Figure 1. Node’s load specification

A P2P system is defined to be balanced if the sum of the load L_i of a physical node P_i is smaller than or equal to the target load of the node for every node $P_i, 1 \leq i \leq M$ in the system. When the system is imbalanced, the goal of a load balancing algorithm is to move VSs from overloaded node to underloaded one with minimum load transfer overheads.

The amount of overload to be transferred from the overloaded node $P_i; 1 \leq i \leq M$ is a random variable denoted by A is given by

$$A(p_i) = \max(0, Q_i - U) = \begin{cases} Q_i - U & \text{if } Q_i > U \\ 0 & \text{Otherwise} \end{cases}$$

Similarly, the amount of underload that can be accepted at the underloaded peer $P_i; 1 \leq i \leq M$ is a random variable denoted by B is given by

$$B(p_i) = \max(0, L - Q_i) = \begin{cases} L - Q_i & \text{if } Q_i < L \\ 0 & \text{Otherwise} \end{cases}$$

Definition 2.

Let $\{Q(t); t \geq 0\}$ be a stationary² stochastic process with state space consisting of non-negative integers. Let S_i and S_j be two distinct non-negative numbers. The First Passage Time (FPT) between states S_i and S_j is denoted by $FPT(S_i, S_j)$, is given by

$$FPT(S_i, S_j) = \begin{cases} \inf \{t; Q(t) = S_j, Q(0) = S_i\} & \text{if } S_i \neq S_j \\ 0 & \text{if } S_i = S_j \end{cases}$$

It is a random variable which measures the minimum amount of time needed to reach state S_j from state S_i . We note that because the same stochastic process $Q(t)$ is stationary, translating the above events by a fixed amount of time has no effect upon the probability distribution of $FPT(S_i, S_j)$. In fact, a first passage time from state S_i to state S_j can be divided into two parts, namely the first transition out of state S_i (say S_k) followed by the first passage from S_k to S_j .

Assume that, $i < j$; since changes of state have unit magnitude in a birth and death of load, then

$$FPT_{ij} = FPT_{ik} + FPT_{kj} \quad i < k < j; k = i+1, i+2, \dots, j-1 \text{ thus} \\ FPT_{ij} = \sum_{k=i}^{j-1} FPT_{k,k+1} \quad i < j \tag{1}$$

Similarly, if $i > j$

$$FPT_{ij} = \sum_{k=j}^{i-1} FPT_{k+1,k} \quad i > j \tag{2}$$

² A stationary stochastic process has the property that the joint distribution don not depend on the time origin. The stochastic process $\{Q(t); t \in \mathcal{T}\}$ is called stationary if $t_i \in \mathcal{T}$ and $t_i + s \in \mathcal{T}, i=1, 2, \dots, k$ (k is any positive integer), then $\{Q(t_1), \dots, Q(t_k)\}$ and $\{Q(t_1+s), \dots, Q(t_k+s)\}$ have the same joint distribution [22].

Let $H_{ij}(t) = P\{FPT_{ij} \leq t\}$ and considering that the summands in both equation (1) and (2) are independent. Thus, if we apply Laplace transformer

$\tilde{H}_{ij} = \prod_{k=i}^{j-1} \tilde{H}_{k,k+1}(s)$ for $H_{ij}(t)$, then we can show that:

$$\tilde{H}_{ij}(s) = \prod_{k=i}^{j-1} \tilde{H}(s, k, k+1) \quad ; i < j \quad (\text{Upward}) \quad (3)$$

$$\tilde{H}_{ij}(s) = \prod_{k=j}^{i-1} \tilde{H}(s, k+1, k) \quad ; i > j \quad (\text{Downward}) \quad (4)$$

Clearly the distribution of the first passage time of unit downward is independent of starting state while the distribution of the first passage time of unit upward depends upon starting state.

B. Load Sharing Edge

The aim of this section is to study the *FPT* of the transition from normal state to overloaded, overloaded to normal, underloaded to normal, etc. For each transfer pair, *FPT* will be computed to predicate the future behavior of the transfer pair before the load transfer (i.e. virtual server migration) decision is taken.

Definition 3.

Let $[Q(t), R(t)]$ be a transfer pair with $Q(t)=X$ and $R(t)=Y$, where $X>U$ and $Y<L$. The **Load Sharing Edge LSE** between Q and R is a random variable $E(X,Y)$ which is defined as follows:

$$E(X,Y) = \min \{FPT(X,U), FPT(Y,L)\}$$

Where, $FPT(X,U)$ is the first passage time to move from state X to state U and $FPT(Y,L)$ is the first passage time to move from state Y to state L .

LSE is the period of time within which the overloaded node must complete transferring load to the underloaded node before the overloaded node identifies that is unnecessary to transfer load or the underloaded node becomes ineligible to receive a transferred load. Assume the load transfer time is denoted by Δ . It is the time needed to package and send the load (i.e. the least loaded virtual server that will release overload) to sink R . Thus, the load transfer must be initiated only if $LSE > \Delta$. Since *LSE* is a random variable we need to formulate the transfer criterion in terms of probabilities. Assume the probability that *LSE* exceeds Δ is $P\{E(X,Y) > \Delta\}$. Therefore, the load transfer must be initiated if $P\{E(X,Y) > \Delta\}$ is large. These considerations led to the formulation of a class of rules so called *Quantile rules*. The quantile of a probability distribution function is defined as follows:

Definition 4.

Let $F(t), t \geq 0$ be the probability distribution function of a non-negative random variable X . Let $0 < \beta < 1$. The β -quantile of F is a non-negative real number q_β satisfying

$$1 - F(q_\beta) = \beta, \quad P\{X \geq q_\beta\} = \beta.$$

From definition 4, the β -quantile rule for load transfer was introduced.

Definition 5.

Given a transfer pair $[Q(t), R(t)]$ and a load transfer time is Δ . Also, for $0 < \beta < 1$, let q_β be the β -quantile of the probability distribution of the *LSE* between $Q(t)$ and $R(t)$. Then the load transfer is initiated only if $q_\beta > \Delta$.

The proposed algorithm in this paper uses a β -quantile rule before transferring load and ensures that $1 \geq P\{E(X,Y) \geq \Delta\} \geq \beta$. In general β can be taken 0.90 or large.

The probability distribution of the random value *LSE* is given as follows

$$P_e(t; X, Y) = P\{E(X, Y) \leq t\} \quad \text{for } t \geq 0.$$

Thus, the probability distribution function of the load sharing edge *LSE* between pair $Q(t)$ and $R(t)$ is given by

$$P_e(t; X, Y) = 1 - [1 - F(t; X, U)] \times [1 - G(t; Y, L)].$$

Where $F(t; X, U)$ and $G(t; Y, L)$ are the probability distribution of the first passage time from X to U and from Y to L in the queues $\{Q(t)\}, \{R(t)\}$, respectively. Each node is modeled as *M/M/1* queue, in which processes arrive according to a Poisson process with mean arrival rate λ , then processed with exponential service time [23, 24] and with mean service rate μ .

Lemma 1.

Assume constant birth rates $\lambda = \lambda_0 = \lambda_1 = \dots$ and death rates $\mu = \mu_0 = \mu_1 = \dots$, then the probability distribution function is

$$H_{k+1,k}(\cdot) = H_{1,0}(\cdot) \quad k=0,1,\dots$$

Proof:

$FPT_{k+1,k}$ is the time that elapses before the cumulative number of deaths first exceeds the cumulative number of births when $X(0)=k+1$. Also, the value of $H_{k+1,k}(\cdot)$ do not depend on $X(0)$, [23, 24]. □

Lemma 2.

The Laplace-Stieljes transform of the probability distribution function of the first passage time from state k to state 0 in an *M/M/1* queue is

$$\tilde{H}_{k,0}(s) = \left[\frac{s + \lambda + \mu - \sqrt{(s + \lambda + \mu)^2 - 4\lambda\mu}}{2\lambda} \right]^k \quad (5)$$

Proof:

Assume that, the first passage time FPT_{ij} can be expressed as

$$FPT_{ij} = S_1 + \begin{cases} FPT_{i+1,j} & \text{if } X(s_1) = i + 1 \\ FPT_{i-1,j} & \text{if } X(s_1) = i - 1 \end{cases}$$

Where S_1 is the time of the first transition. Assume that $H_{ij}(t) = P\{FPT_{ij} \leq t\}$, as well FPT_{ij} can be upward or downward after the first transition S_1 . Using theorem 4-7, [23, 24], $H_{ij}(t)$ can be expressed as follows

$$H_{ij}(t) = \lambda_i \int_0^t H_{i+1,j}(t-x) e^{-(\lambda_i+\mu)x} dx + \mu_i \int_0^t H_{i-1,j}(t-x) e^{-(\lambda_i+\mu)x} dx \quad (i)$$

Taking Laplace-Stieljes transform on both side of equation (i) and use the convolution property, then the following equation can be obtained:

$$\tilde{H}_{i,j}(s) = \frac{\lambda_i \tilde{H}_{i+1,j}(s) + \mu_i \tilde{H}_{i-1,j}(s)}{s + \lambda_j + \mu_j} \quad (ii)$$

Thus,

Set $i=1$, and $j=0$ $\tilde{H}_{1,0}(s) = \frac{\lambda_0 \tilde{H}_{2,0}(s) + \mu_0 \tilde{H}_{0,0}(s)}{s + \lambda_0 + \mu_0}$ (iii)

From lemma 1, $\lambda = \lambda_0$, $\mu = \mu_0$ and from equation (4)

$\tilde{H}_{i,j}(s) = \prod_{k=j}^{i-1} \tilde{H}_{k+1,k}(s)$; $i > j$. Thus,

$\tilde{H}_{2,0}(s) = \tilde{H}_{1,0}(s) \times \tilde{H}_{2,1}(s) = (\tilde{H}_{1,0}(s))^2$ from equation (iii), we obtain the following quadratic equation

$\lambda [\tilde{H}_{1,0}(s)]^2 - (s + \lambda + \mu) \tilde{H}_{1,0}(s) + \mu = 0$ (iv)

Equation (iv) has two solutions, we consider the solution which satisfies that $\tilde{H}_{1,0}(s) \leq 1$

$\tilde{H}_{1,0}(s) = \left[\frac{s + \lambda + \mu - \sqrt{(s + \lambda + \mu)^2 - 4\lambda\mu}}{2\lambda} \right]; \tilde{H}_{1,0}(s) \leq 1$, for real.

Hence, $\tilde{H}_{k,0}(s) = [\tilde{H}_{1,0}(s)]^k = \left[\frac{s + \lambda + \mu - \sqrt{(s + \lambda + \mu)^2 - 4\lambda\mu}}{2\lambda} \right]^k$

□

Corollary 1.

The density function of the first passage time from state k to state 0 in an M/M/1 queue is

$h_{k,0}(t) = ke^{-(\lambda+\mu)t} I_k(2t\sqrt{\lambda\mu}) \frac{(\mu/\lambda)^{k/2}}{t}$ (6)

For $t > 0$ m; where I_k is the Modified Bessel function of order k .

Proof:

If equation (5) is bona fide Laplace transform, it is the Laplace transform of $h(t)$, [23, and 24]. From Laplace transform,

$\wp\{e^{-ct} f(t)\} = f(s+c)$; set $(c = \lambda + \mu)$ and $w = (s + \lambda + \mu)$;

So we can write equation (5) as

$\tilde{H}_{1,0}(w) = \left[\frac{w - \sqrt{w^2 - 4\lambda\mu}}{2\lambda} \right]$, [23, 24].

While the Laplace transformation

$\wp\{I_1(at)/t\} = \frac{s + \sqrt{s^2 + a^2}}{a}$, [23] then we can say that the

numerator is the Laplace transform of $2(\lambda\mu)^{1/2} I_1(2t(\lambda\mu)^{1/2})/t$; thus, $\tilde{H}_{1,0}(s)$ is the Laplace transform of equation (5). So we have

$h_{1,0}(t) = ke^{-(\lambda+\mu)t} I_1(2t\sqrt{\lambda\mu}) \frac{\sqrt{(\mu/\lambda)}}{t}$.

Hence,

$h_{k,0}(t) = ke^{-(\lambda+\mu)t} I_k(2t\sqrt{\lambda\mu}) \frac{(\mu/\lambda)^{k/2}}{t}$; $t > 0$. □

Lemma 3. (Downward)

The probability distribution of the first passage time from state k to state 0 in an M/M/1 queue is:

$H(x; k, 0) = 1.0 - k\mu^k \sum_{n=0}^{\infty} \frac{(\lambda\mu)^n}{(\lambda + \mu)^{2n+k}} \Gamma(2n + k, x)$

(7)

Where $\Gamma(\)$ is the incomplete Gama function.

Proof:

Assume that the modified Bessel function of order k is

$I_k(x) = \left(\frac{x}{2}\right)^k \sum_{n=0}^{\infty} \frac{(x^2/4)^n}{n!(n+k)!}$; $x \geq 0$

Set $x = 2t\sqrt{\lambda\mu}$. Thus,

$I_k(2t\sqrt{\lambda\mu}) = \left(\frac{2t\sqrt{\lambda\mu}}{2}\right)^k \sum_{n=0}^{\infty} \frac{(t^2\lambda\mu)^n}{n!(n+k)!}$; $x \geq 0$ By substitute

into equation (6), we can compute the density function $h_{k,0}(t)$ as follows:

$h_{k,0}(t) = k\mu e^{-(\lambda+\mu)t} \sum_{n=0}^{\infty} \frac{(\lambda\mu)^n t^{k+2n-1}}{n!(n+k)!}$; $t \geq 0$

From the definition of the probability distribution function, we have $H_{k,0}(x) = \int_0^x h_{k,0}(t) dt$ and

$H_{k,0}(x) = 1 - \int_x^{\infty} h_{k,0}(t) dt$.

Thus, $H_{k,0}(x) = 1 - \int_x^{\infty} k\mu^k e^{-(\lambda+\mu)t} \sum_{n=0}^{\infty} \frac{(\lambda\mu)^n t^{k+2n-1}}{n!(n+k)!} dt$.

If we exchange the infinite sum and the integral we get

$H_{k,0}(x) = 1 - \sum_{n=0}^{\infty} \frac{(\lambda\mu)^n}{n!(n+k)!} \int_x^{\infty} e^{-(\lambda+\mu)t} t^{k+2n-1} dt$; $\Gamma(i, x) = \int_x^{\infty} e^{-t} t^{i-1} dt$.

Hence,

$H(x; k, 0) = 1.0 - k\mu^k \sum_{n=0}^{\infty} \frac{(\lambda\mu)^n}{(\lambda + \mu)^{2n+k}} \Gamma(2n + k, x)$. □

Lemma 4.

The probability density function of the first passage time of the M/M/1 queue from state 0 to state 1 is $h_{0,1}(t) = \lambda e^{-\lambda t}$.

Proof:

Let $\tilde{h}_{i,j}(s) = \frac{\lambda_i \tilde{h}_{i+1,j}(s) + \mu_i \tilde{h}_{i-1,j}(s)}{s + \lambda_j + \mu_j}$ set $i=0, j=1$ (i.e. there

is no service $\mu_0=0$) and $\tilde{h}_{0,1}(s) = \frac{\lambda_0 \tilde{h}_{1,1}(s)}{s + \lambda_0}$. For

simplicity set $\lambda_0 = \lambda$. Hence $\tilde{h}_{0,1}(s) = \frac{\lambda}{s + \lambda}$; $\tilde{h}_{1,1} = 1$ since

$FPT_{11}=0$. By computing the inverse of the Laplace transformation we get:

$h_{0,1}(t) = \lambda e^{-\lambda t}$. □

Lemma 5. (Upward)

The probability distribution function of the first passage time of the M/M/1 queue from state i to state j ; $i < j$ is

$H_{i,j}(t) = 1 + \lambda^{j-i} \sum_{k=1}^j C_k e^{-r_k t}$

(8)

where r_k ; $1 \leq k \leq j$, are j distinct roots of the polynomial of degree j defined recursively as:

$D(s-1) = 0, D(s,0) = 1; D(s,1) = s + \lambda$

$D(s, j) = (s + \lambda + \mu)D(s, j-1) - \lambda\mu D(s, j-2); j \geq 2$

also, for $1 \leq k \leq j$, $C_k = (s + r_k) \frac{D(s, i)}{sD(s, i)} \Big|_{s=-r_k}$.

Proof:

In lemma 4, we have prove $\tilde{h}_{0,1}(s) = \frac{\lambda}{s + \lambda}$ that is the

Laplace transformation of the probability distribution of the first passage time of FPT_{01} . Assume that

$\tilde{h}_{i,j}(s) = \frac{\lambda_i \tilde{h}_{i+1,j}(s) + \mu_i \tilde{h}_{i-1,j}(s)}{s + \lambda_j + \mu_j}$. Set $i = k, j = k + 1$,

$\lambda_k = \lambda_{k+1} = \lambda, \mu_k = \mu_{k+1} = \mu$

$$\tilde{h}_{k,k+1}(s) = \frac{\lambda \tilde{h}_{k+1,k+1}(s) + \mu \tilde{h}_{k-1,k+1}(s)}{s + \lambda + \mu} \tag{i),}$$

But, $\tilde{h}_{k+1,k+1}(s) = 1$ since $FPT_{k+1, K+1} = 0$,

$\tilde{h}_{k-1,k+1}(s) = \tilde{h}_{k-1,k}(s) \tilde{h}_{k,k+1}(s)$. Hence,

$$\tilde{h}_{k,k+1}(s) = \frac{\lambda + \mu \tilde{h}_{k-1,k}(s) \tilde{h}_{k,k+1}(s)}{s + \lambda + \mu}$$

$$\tilde{h}_{k,k+1}(s) = \frac{\lambda}{s + \lambda + [\mu(1 - \tilde{h}_{k-1,k}(s))]} \tag{ii),}$$

Using mathematical induction, we can prove that equation (ii) is satisfied for all $k \geq 0$. It can be rewritten as the ratio of two functions $N(s, k)$ and $D(s, k)$.

These functions can be defined as follows:

$$L(1) = -\lambda, \quad L(k) = -(\lambda + \mu), k \geq 2$$

$$M(k) = \lambda\mu \quad k \geq 1$$

Thus,

$$N(s, k) = \lambda D(s, k - 1), k > 1; D(s, 0) = 1$$

$$D(s, k) = [s - L(k)] \times D(s, k - 1) - M(k) \times D(s, k - 2); k \geq 1$$

Where $D(s, k)$ is a polynomial of degree $k; k \geq j$.

Thus, equation (ii) can be rewritten as follows:

$$\tilde{h}_{k-1,k}(s) = \frac{N(s, k)}{D(s, k)} = \frac{\lambda D(s, k - 1)}{s + \lambda + \mu D(s, k - 1) - \lambda \mu D(s, k - 2)}, k \geq 2$$

But for general transitions from state i to state j : $0 \leq i \leq j$, is

$$\tilde{h}_{i,j}(s) = \frac{\lambda D(s, i)}{D(s, i + 1)} \times \frac{\lambda D(s, i + 1)}{D(s, i + 2)} \times \dots \times \frac{\lambda D(s, j - 2)}{D(s, j - 1)} \times \frac{\lambda D(s, j - 1)}{D(s, j)}$$

Thus, we can obtain the Laplace transform of the density of the upward transition from state i to state j ,

$$\tilde{h}_{i,j}(s) = \frac{\lambda^{j-i} D(s, i)}{s D(s, j)}$$

the denominator and the numerator from the above equation.

From the definition $LT \left[\int_0^x h(t) dt \right] (s) = \frac{LT[h](s)}{s}$, thus

$$\tilde{H}_{i,j}(s) = \frac{\lambda^{j-i} D(s, i)}{s D(s, j)}$$

It is a relation function in which the numerator polynomial has degree I while the denominator polynomial has degree $(j+I)$ and $(i < j)$. Thus, we can expand $\tilde{H}_{i,j}(s)$ into a finite sum of partial fractions as follows:

If $D(0, j) = \lambda^j$ for all $j \geq 1$ and $\lambda > 0$ then zero cannot be a root of $D(s, j); j \geq 1$ and it can be written in the following form:

$$D(s, j) = (s + r_1)(s + r_2) \dots (s + r_j) \text{ where } r_k \text{ for all } k = 1$$

to j are roots of $D(s, j)$. Hence, it can be written as

$$\tilde{H}_{i,j}(s) = \lambda^{j-i} \left(\frac{C_0}{s} + \sum_{k=1}^j \frac{C_k}{(s + r_k)} \right)$$

follows:

Where $C_k = (s + r_k) \frac{D(s, i - 1)}{s D(s, j)} \Big|_{s = -r_k}; 1 \leq j \leq k$ and

$$C_0 = \frac{D(0, i)}{D(0, j)} = \frac{\lambda^i}{\lambda^j} = \lambda^{i-j}$$

$$\tilde{H}_{i,j}(s) = \frac{1}{s} + \lambda^{j-i} \sum_{k=1}^j \frac{C_k}{(s + r_k)}$$

Accordingly, we can invert the Laplace Transform of the above equation. But, each term in the right hand side is in the form $\frac{\alpha}{(s + \beta)^{\alpha, \beta}}$ Where α, β are constants.

Each term has inverse Laplace transformation $\alpha e^{-\beta t}$. Hence,

$$H_{i,j}(t) = 1 + \lambda^{j-i} \sum_{k=1}^j C_k e^{(-r_k t)} \quad \square$$

Theorem

Let $[Q(t), R(t)]$ be a transfer pair that consists of M/M/1 queues. Let m be the amount of overload and n the amount of underload. Then the probability distribution function of the Load Sharing Edge (LSE) is

$$P_e(t; m, n) = 1 - m \lambda^n \mu^m \left[\sum_{k=1}^j C_k e^{-r_k t} \right] \times \left[\sum_{k=0}^{\infty} \frac{(\lambda \mu)^k}{(\lambda + \mu)^{2k+m}} \Gamma(2n + m, x) \right]$$

where, $r_k, 1 \leq k \leq L$; are the roots of the polynomial defined recursively as

$$D(s, -1) = 0; D(s, 0) = 1; D(s, 1) = s + \lambda$$

$$D(s, L) = (s + \lambda + \mu) D(s, L - 1) - \lambda \mu D(s, L - 2); L \geq 2; 1 \leq k \leq L$$

Also, for $1 \leq k \leq L$ $C_k = (s + r_k) \frac{D(s, L - n)}{s D(s, L)} \Big|_{s = -r_k}$.

Proof:

The transfer pair has the probability distribution function $P_e(t; i, j)$ of the LSE; $i > U, j < L$ which is defined by

$P_e(t; i, j) = 1 - [1 - F(t; i, U)] * [1 - G(t; j, L)]$ where m is the amount of overload ($m = i - U$), and n is the amount of underload ($n = L - j$) then

$P_e(t; m, n) = 1 - [1 - F(t; m, 0)] * [1 - G(t; L - n, L)]$ for an M/M/1 queue case $F(t, m, 0)$ and $G(t, L - n, L)$ have been derived from 3 and 5 respectively. Hence,

$$P_e(t; m, n) = 1 - m \lambda^n \mu^m \left[\sum_{k=1}^L C_k e^{-r_k t} \right] \times \left[\sum_{k=0}^{\infty} \frac{(\lambda \mu)^k}{(\lambda + \mu)^{2k+m}} \Gamma(2n + m, x) \right]$$

□

Due to the infinite number of terms in the probability distribution $P_e(t; m, n)$ of the LSE in M/M/1, the following lemma will drive a formula for LSE as finite number of terms as follows.

Lemma 6

For a transfer pair $[Q(t), R(t)]$ with an amount of overload m and an amount of underload n , the Mean Load Sharing Edge MLSE(m, n) is

$$MLSE(m, n) = -\lambda^n \sum_{k=1}^L C_k \left(\frac{1}{r_k} - \tilde{F}(r_k, m, 0) \right)$$

where $r_1, r_2, \dots, r_k, C_k$ are constants defined in the previous theorem. Also,

$$\tilde{F}_{m,0}(s) = \left[\frac{s + \lambda + \mu - \sqrt{(s + \lambda + \mu)^2 - 4\lambda\mu}}{2\lambda} \right]^m$$

Proof:

Since $P_e(t; m, n) = 1 - [1 - F(t; m, 0)] * [1 - G(t; L - n, L)]$ we obtain,

$$MLSE(m, n) = \int_0^\infty P_e(t, m, n) dt$$

$$= \int_0^\infty [1 - F(t; m, 0)] \times [1 - G(t; L - n, L)] dt$$

From lemma 5, we get

$$MLSE(m, n) = - \int_0^\infty \left[\lambda^n \sum_{k=1}^L C_k e^{-r_k t} \right] \times [1 - F(t; m, 0)] dt$$

But the Laplace transform $LT[q](s)|_{s=0} = \int_0^\infty q(t) dt$ and

$$LT[e^{-\alpha t} q](s)|_{s=0} = LT[q](s + \alpha). \text{ Set } \alpha = r_k, \text{ then}$$

$$MLSE(m, n) = - \lambda^n \sum_{k=1}^L C_k \left(\frac{1}{r_k} - \tilde{F}(r_k, m, 0) \right) \quad \square$$

Consequently, it has been observed that for queuing models in which job arrival and processing rates are independent of queue size, such as M/M/1 queues. The distribution of LSE depends only on the amounts of underload and overload. The following algorithm will use the numeric value of the given mean load sharing edge formula that is based upon the following parameters $(m, n, \lambda_i, \mu_i, \lambda_j, \mu_j, L, q_\beta)$.

C. Algorithm

This section introduces the proposed load balancing algorithm based upon the above analysis. Periodically every T seconds, each overloaded peer transfers the exceeds load to the underloaded peers (i.e. sender-initiated algorithm). This algorithm imposes a β -quantile rule for transferring load. For each pair (λ, μ) a corresponding β -quantile should be determined while β must be taken 0.90 or large. The algorithm is shown in the following scenario:

1. The overloaded node S_i creates a suitable domain (group) D_i from neighbor nodes to the peer S_i . Each node belongs to D_i satisfies $D_i = \{S_j; P(FPT_i > t_{ij}) \geq 0.90 \text{ and } i \neq j\}$. Where t_{ij} is the required communication delay to send a message from node S_i to S_j plus the required time to reply with load transfer of certain virtual server from S_i to S_j . Also, FPT_i is the first passage time of node S_i to transfer from overloaded state to normal or underloaded state. D_i is an ordered set with respect to the communication t_{ij} . It is implemented as an order linked list.
2. Thus, S_i sends a broadcast messages to all nodes belonging to the domain D_i . Node S_i must receive a reply from all nodes belonging to D_i within the FPT_i time.
3. Node S_i selects an underloaded node $S_j \in D_i$ where the mean load sharing edge $MLSE$ between S_i and S_j . if $q_\beta > \Delta$ then transfers load (virtual server) from S_i to S_j . Where pair λ and μ are given, Δ is the time needed to transfer load less than or equal to $A(S_j)$, β is 0.90 or large.
4. Repeat step 3 for each underloaded node S_j belonging to D_i whenever FPT_i period doesn't run out yet.

```

Load_balance( $S_i, T$ )
{
  // Every period T seconds  $S_i$  checks its load
  // jumps above upper load U.
  // It should do the following.
  Create( $D_i, S_i$ ); // create domain of neighbors
  While ( $q_\beta > 0$ ) do
    // repeat the following within a period  $q_\beta$ 
    Select  $S_j \in D_i$ ; // select from order set  $D_i$ 
     $D_i = D_i \setminus \{S_j\}$ 
    If ( $q_\beta > \Delta$ ) then transfer_load( $S_i, S_j$ );
  }
}

transfer_load( $S_i, S_j$ )
{
  If !(Overloaded) then return; //Sender-initiated
  If ( $S_i \rightarrow VSset.size > 1$ ) then
    Choose  $v \in S_i \rightarrow VSset$  such that:
    a. Transfer  $v$  to  $S_j$ ; will not overload  $S_j$ 
    b.  $v$  is the least loaded virtual server that will release
       overload.
    Failing that, let  $v$  be the most loaded VS.
    Return the virtual servers reassignment.
}
    
```

IV. EVALUATION

This paper implements an event-based simulation to evaluate the proposed load balancing algorithm. It uses several parameters as follows: default number of virtual servers per node (12), number of nodes (4096), system

utilization (0.8), Object arrival rate (Poisson with mean arrival time 0.01 sec), average number of objects (1 million), and periodic load balancing period (T=60 seconds). This simulation evaluates the following metric. Load Movement Ratio (LMR), defined as the total movement cost incurred due to load balancing divided by the total cost of moving all objects in the system at once. In case the value of the LMR is 0.1, it infers that the balancer consumes about 10% of its bandwidth to insert objects. The node arrival rate is modeled by a Poisson process, and the lifetime of a node is drawn from an exponential distribution. This simulation ran with two inter-arrival times 10 and 60 seconds, it fixes the steady-state number of nodes in the system to 4096 nodes.

Fig. 2 plots the LMR metric as a function of system utilization, to study the load moved by the proposed load balancing algorithm as a fraction of the load moved by the underlying DHT due to node arrivals and departures. Fig. 2 demonstrates that the load moved by the proposed load balancing algorithm is significantly smaller than the load moved by the underlying DHT especially for small system utilization. In addition, fig. 2 shows that the LMR with node inter-arrival time 10 sec is larger than with node inter-arrival time 60 sec. Fig. 3 verifies the perception that increasing the number of virtual servers decreases considerably the fraction of load moved by the underlying DHT. Fig. 4 demonstrates that increasing

number of virtual servers per node assists load balance at high system utilizations and grants efficient load movements due to low bandwidth lost. Fig. 5 plots the 90-quantile of the load sharing edge (LSE) with system utilization when overload is 5 at source node and underload is 2 at the destination node. It demonstrates that the 90-quantile of the LSE tends to be smaller as system utilization increases. As seen from fig. 5, the 90-quantile of the LSE is 9.94ms thus load can be transferred only if $\Delta < 9.94\text{ms}$. In addition, increasing the number of virtual servers reduces significantly the 90-quantile that helps in avoiding load thrashing.

V. CONCLUSION

Load balancing among peers is critical and a key challenge in peer-to-peer systems. This paper exhibits a stochastic analysis that avoids the load thrashing and tackles the out-of-date problems due to peer's state changes during load movement (virtual servers migration). Then, it proposes a load balancing algorithm based on that stochastic analysis. An efficient simulation has been carried out that demonstrates the effectiveness of the proposed load balancing algorithm.

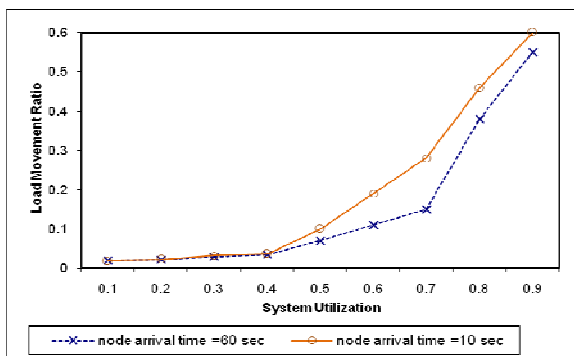


Figure 2. LMR vs. system utilization with two node arrival times.

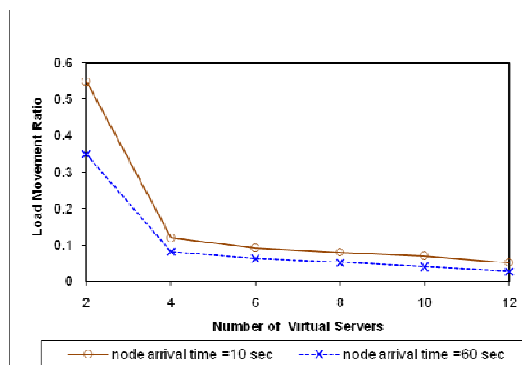


Figure 3. LMR vs. number of virtual servers with two node arrival times.

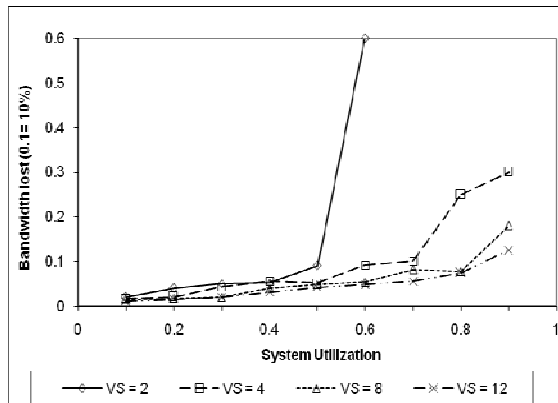


Figure 4. Bandwidth lost vs. system utilization with different number of virtual servers per node.

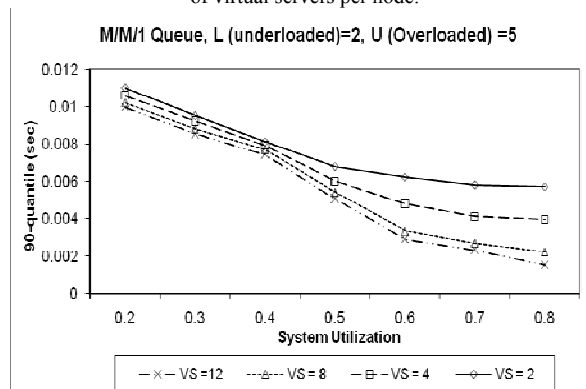


Figure 5. 90-quantile of the LSE vs. system utilization with different number of virtual server per node.

REFERENCES

- [1] Shirky C., "Modern P2P Definition," <http://www.openp2p.com/pub/a/p2p/2000/11/24/shirky1-whatisp2p.html>, 2000.
- [2] S. Saroiu, et. al., "A Measurement study of Peer-to-Peer sharing systems," Proc. Multimedia Computing and Networking Conf. (MMCN), 2002.
- [3] I. Sotica, et. al., "Chord: A scalable peer-to-peer lookup service for Internet Applications," ACM SIGCOMM'01, August 2001, pp. 149-160.
- [4] Antony Rowstron and Peter Druschel, "Pastry: Scalable distributed object location and routing for large-scale Peer-to-Peer systems," in Proc. Middleware, 2001.
- [5] S. Ratnasamy, et al. "A Scalable Content- Addressable Network," in Proc. ACM SIGCOMM'01, California, USA
- [6] Derk I., et al. "Adaptive load sharing in homogenous distributed systems," IEEE Trans. on Soft. Eng., Vol. 12 No 5, Sept. 1986.
- [7] A. R. Bhambe, et. al., "Mercury: supporting scalable multi-attribute range queries," In Proc. of the Conf. on Applications, Technologies, Architectures, and Protocols For Computer Communication, ACM, New York, 2004.
- [8] Frank Dabek, et. al., "Wide-area cooperative storage with CFS," Proc. 18th ACM Symp. Operating Systems Principles (SOSP'01), pp. 2020-215, Oct. 2001
- [9] Godfrey et. al., "Load balancing in dynamic structured P2P systems," Proc. IEEE INFOCOM, 2004.
- [10] J. Byers, "Simple load balancing for distributed hash table," Proc. Of the 2nd Int. Workshop on Peer-to-peer systems (IPTPS'03), Feb. 2003.

- [11] A. Rao et. al., "Load balancing in structured P2P systems," Proc. Of the 2nd Int. Workshop on Peer-to-peer systems (IPTPS'03), Feb. 2003.
- [12] P.B.Godfrey and I. Stoica, "Heterogeneity and Load balance in distributed hash table," Proc. IEEE INFOCOM, 2005.
- [13] Jie Li, Hisao Kameda, "A decomposition algorithm for optimal static load balancing in Tree hierarchy network configurations," IEEE Trans. On Parallel and Distributed Systems, Vol. 5, No.5, 1994.
- [14] T.L. Casavant, and J.G. Kuhl, "A taxonomy of scheduling in general-purpose distributed computing systems," IEEE Trans. on Software Engineering, vol. 14, no. 2, pp. 141-154, Feb., 1988.
- [15] N. G. Shivaratri, et. al. , "Load distributing for locally distributed systems," Computer, vol.25, no.12, pp.33-44, Dec 1992
- [16] Goscinski, A., "Distributed Operating System: The logical design," Addison-Wesly, 1991.
- [17] S. Zhou, "A Trace-Driven Simulation Study of Dynamic Load Balancing," IEEE Transactions on Software Engineering, vol. 14, no. 9, pp. 1327-1341, Sept., 1988.
- [18] D. L. Eager, et.al., "A comparison of receiver-initiated and sender-initiated adaptive load sharing," SIGMETRICS Performance Evaluation Review 13, 2, Aug. 1985.
- [19] P. Kruger, and N. G. Shivaratri, "Adaptive location policies for global schedule," IEEE Soft. Eng., Vol 20, 432-443, June 1994.
- [20] S. Zhou, et.al., "Utopia: a load sharing facility for large, heterogeneous distributed computer systems," Software-Practice and Experience. Dec. 1993, pp. 1305-1336.
- [21] S. P. Dandamudi, and K. C. Lo, "A Hierarchical Load Sharing Policy for Distributed Systems," in Proc. of the 5th Int. Workshop on Modeling, Analysis, and Simulation of Computer and Telecommunications Systems, MASCOTS. IEEE CS, Washington, DC, 1997.
- [22] Heyman, Daniel P., "Stochastic models in operation research," vol. I, MxGraw-Hill Inc. 1982.
- [23] Hisashi Kobayashi, "Modeling and Analysis: An Introduction to system performance evaluation methodology," Addison-Wedley, 1978.
- [24] Hisashi Kobayashi, et. al., "System modeling and analysis: foundations of system performance evaluation," Prentice Hall 2009.
- [25] Zhu Y. and Hu Y., "Efficient, Proximity-aware Load Balancing for DHT-based P2P Systems," IEEE Trans. On Parallel and Distributed Systems, Vol. 16, No. 4, 2005



Khaled Ragab is an assistant professor at Department of Computer Science, College of Computer and Information Technology, King Faisal University, Saudi Arabia. Moreover, he is on leave assistant professor of Computer Science at Department of Mathematic, Computer Science division, Ain Shams University.

He joined Department of Computer Science, Tokyo University in 2005 as postdoctoral position. He was born in 1968 and received his B.Sc., M.Sc. degrees in Computer Science from Ain Shams University, Cairo, Egypt in 1990, 1999, respectively and Ph.D. degree in Computer Science from Tokyo Institute of Technology in 2004. He has worked in Ain Shams University, Cairo Egypt in 1990-1999 as assistant lecturer. He has worked as research scientist in Computer Science Dept., Technical University of Chemnitz, Germany in 1999-2001. His research interests include autonomous decentralized systems, Peer-to-Peer Systems, Overlay Networks, Load Balancing, Web-services and application-level multicast.