

Proximity-Aware Superpeer Overlay Topologies

Gian Paolo Jesi, Alberto Montresor, and Ozalp Babaoglu

Abstract—The concept of *superpeer* has been introduced to improve the performance of popular P2P applications. A *superpeer* is a “powerful” node that acts as a server for a set of clients, and as an equal with respect to other superpeers. By exploiting heterogeneity, the superpeer paradigm can lead to improved efficiency, without compromising the decentralized nature of P2P networks. The main issues in constructing superpeer-based overlays are the selection of superpeers and the association between superpeers and clients. Generally, superpeers are either run voluntarily (without an explicit selection process), or chosen among the “best” nodes in the network, for example those with the most abundant resources, such as bandwidth or storage. In several contexts, however, shared resources are not the only factor; latency between clients and superpeers may play an important role, for example in *online games* and *IP-Telephony applications*. This paper presents SG-2, a novel protocol for building and maintaining proximity-aware superpeer topologies. SG-2 uses a gossip-based protocol to spread messages to nearby nodes and a biology-inspired task allocation mechanism to promote the “best” nodes to superpeer status. The paper includes extensive simulation experiments to prove the efficiency, scalability and robustness of SG-2.

Index Terms—P2P, superpeer, overlay, latency, quality of service.

I. INTRODUCTION

MODERN P2P networks present several unique aspects that distinguish them from traditional distributed systems. Networks comprising hundreds of thousand of peers are not uncommon. A consequence of such scale is extreme dynamism, with a continuous flow of nodes joining or leaving. Such characteristics present several challenges to the developer. Neither a central authority nor a fixed communication topology can be employed to control the various components. Instead, a dynamically changing overlay topology is maintained and control is completely decentralized. The topology (i.e., the relation “who knows whom” among nodes) is defined by “cooperation” links among nodes, that are created and deleted based on the requirements of the particular application.

The choice of a particular topology is a crucial aspect of P2P design. Until recently, most deployed P2P applications were characterized by the absence of a specific mechanism for enforcing a given topology; for example, Gnutella nodes were free to accept/refuse connections at will [1]. Inefficient communication schemes, such as flooding, are a consequence of this choice.

A distinct, but related problem regards roles that nodes may assume: original P2P systems were based on a complete

“democracy” among nodes: “everyone is a peer”. But physical hosts running P2P software are usually very heterogeneous in terms of computing, storage and communication resources, ranging from high-end servers to low-end desktop machines.

The superpeer paradigm is an answer to both issues [1], [2]. It is based on a two-level hierarchy: *superpeers* are nodes faster and/or more reliable than “normal” nodes that take on server-like responsibilities and provide services to a set of *clients*. For example, in the case of file sharing, a superpeer builds an index of the files shared by its clients and participates in the search protocol on their behalf. Clients are leveraged from taking part in costly protocols and the overall traffic is reduced by forwarding queries only among superpeers. Superpeers allow decentralized networks to run more efficiently by exploiting heterogeneity and distributing load to machines that can handle the burden. On the other hand, this architecture does not inherit the flaws of the client-server model, as it allows multiple, separate points of failure, increasing the health of the P2P network.

The superpeer paradigm is not limited to file sharing: it can be seen as a general approach for P2P networking. Yet, the structural details are strongly application-dependent, so we cannot identify a “standard” superpeer topology. Parameters to be considered include: how superpeers are linked together; how to arrange clients; how many superpeers are needed; etc.

In this paper, we focus our investigation on a specific aspect of the problem: *proximity*. Our goal is to build a topology where clients and superpeers are related based on their distance (in terms of communication latency). The idea is to select superpeers among the most powerful nodes, and to associate them with clients whose round-trip time is bounded by a specified constant. This is a generic problem, whose solution can be beneficial to several P2P applications. Examples include online games such as Age of Empires [3], P2P telephony networks such as Skype [4] and streaming applications such as PeerCast [5]. In all these cases, communication latency is one of the main concerns.

Our solution, called SG-2, is a self-organizing, decentralized protocol capable of building and maintaining superpeer-based, proximity-aware overlay topologies. SG-2 uses an epidemic protocol to spread messages to nearby nodes, and implements a task allocation protocol that mimics the behavior of social insects. These biology-inspired mechanisms are combined to promote the “best” nodes to the superpeer status, and to associate them to nearby clients.

To validate the results of our protocol, we considered two specific test scenarios: *online games* and *P2P IP-Telephony applications*. In the former, a large number of players interact together (or against each other) in virtual worlds. Most online games follow a client-server model, where the only function

Manuscript received September 1, 2006; revised January 31, 2007; accepted August 1, 2007. The associate editors coordinating the review of this letter and approving it for publication were A. Keller and J. P. Martin-Flatin.

Gian Paolo Jesi and Ozalp Babaoglu are with the University of Bologna (e-mail: jesi@cs.unibo.it).

Alberto Montresor is with the University of Trento.
Digital Object Identifier 10.1109/TNSM.2007.070904.

of the client software is to present a graphic user interface to the player, while the state of the simulated persistent world is hosted on the server side. This approach is scalable only thanks to the deployment of high-end clusters of replicated servers. A small number of games have attempted a different approach. MiMaze [6] and Age of Empires [3] are completely decentralized, and the game state is replicated at all participants. In this case, consistency requirements limit the number of players that may be involved in the same game.

In the P2P Telephony scenario instead, a huge number of users can locate another (known) party (or a small group of users) and start a voice conversation as they would do with normal telephone equipment. The only widely adopted applications of this kind is Skype [7], but it is a closed software and therefore there are not many details about its inner working. However, Skype claims to be a P2P applications based on superpeers. Other solutions are based on the open SIP [8] protocol standard, but this technology is mainly based on the client-server paradigm. However, the idea to merge SIP and P2P is gathering credit [9], [10].

We believe that the superpeer paradigm could represent an interesting alternative to the approaches above. We envision a system where a small number of powerful nodes act as state servers when needed, with the remaining ones acting as clients. All nodes run the same code and can switch from the first role to the second when needed. Thus, superpeers dynamically change over time, depending on the environment conditions.

II. SYSTEM MODEL

We consider a network consisting of a large collection of *nodes*. The network is highly dynamic; new nodes may join at any time, and existing nodes may leave, either voluntarily or by *crashing*. Since voluntary leaves may be simply managed through “logout” protocols, in the following we consider only node crashes. Byzantine failures, with nodes behaving arbitrarily, are excluded from the present discussion. We assume nodes are connected through an existing routed network, such as the Internet, where every node can potentially communicate with every other node. To actually communicate with another node, however, a node must know its *identifier*, e.g. a pair (IP address, port).

The nodes known to a node are called its *neighbors*, and as a set are called its *view*. Together, the views of all nodes define the topology of the overlay network. Given the large scale and the dynamism of our envisioned system, views are typically limited to small subsets of the entire network. Views can change dynamically, and so the overlay topology.

Nodes are heterogeneous: they differ in their computational and storage capabilities, and also (and more importantly) with respect to the bandwidth of their network connection. To discriminate between nodes that may act as superpeers and nodes that must be relegated to the role of clients, each node v is associated with a *capacity* value $cap(v)$, that represents the number of clients that can be handled by v . To simplify our simulations, we assume that each node knows its capacity. In reality, this parameter is strongly dependent on the specific application, and can be easily computed on-the-fly through on-line measurements.

Besides capacity associated to each single node (“how many”), another parameter to be considered is the end-to-end latency between nodes (“how well”). In our model, each pair of nodes (v, w) is associated with a *latency distance* $lat(v, w)$, representing the average round-trip time (RTT) experienced by communications between them. The latency distance between a specific pair of nodes may be measured directly and precisely through ping messages, or approximately estimated through a *virtual coordinate service* [11]; given the dynamic nature of our system and the large number of nodes to be evaluated as potential neighbors, we will adopt the latter approach.

III. THE PROBLEM

Generally speaking, our goal is to create a topology where the most powerful nodes (in terms of capacity) are promoted to the role of superpeers, and the association clients/superpeers is such that each client obtains a configurable *quality of service* (in terms of latency distance) from its superpeer.

More formally, we define the problem of building a proximity-aware, superpeer-based topology as follows. At any given time, the problem input is given by the current set of nodes \mathcal{V} , and the functions $cap()$ and $lat()$ defined over it. Furthermore, a global parameter tol expresses the maximum latency distance that can be tolerated between clients and superpeers. The constraints describing our target topology are the following:

- each node is either a superpeer or a client;
- each client c is associated to exactly one superpeer s (we write $super(c) = s$);
- the number of clients associated to a superpeer s does not exceed $cap(s)$;
- given a superpeer s and one of its clients c , we require that $lat(s, c) \leq tol$.

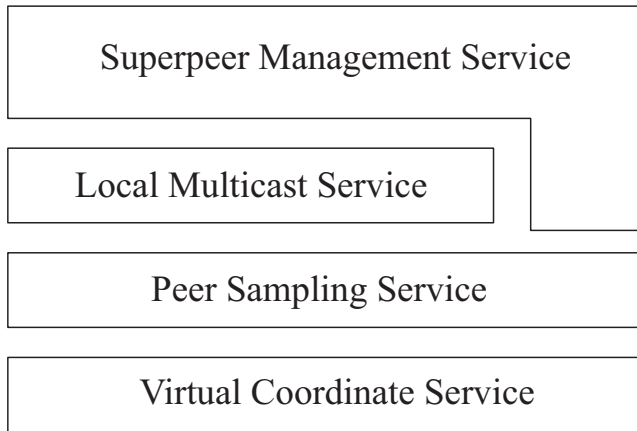
To avoid to end up with a set of disconnected, star-shaped components rooted at each superpeer, we require that superpeers form another proximity-based overlay: two superpeers are connected if their latency distance is smaller than $tol + \delta$, where δ is another configuration parameter.

We aim at selecting as few superpeers as possible (otherwise, the problem could be trivially solved by each node acting as a superpeer, with no client/superpeer connections). This choice is motivated, once again, by the particular scenarios we are considering; for example, in online games, superpeers manage the distributed simulation state, so centralizing as many decisions as possible is important from the performance point of view. Note that given the dynamism of our environment, obtaining the minimum number of superpeers may be difficult, or even impossible. But even in a steady state, the resulting optimization problem is NP-complete.

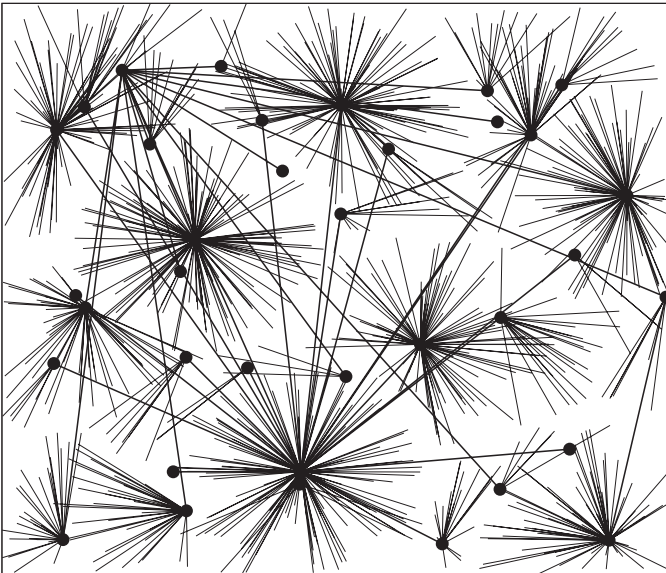
IV. THE SG-2 PROTOCOL

The architecture of SG-2 is shown in Figure 1(a); here, we briefly describe the rationale behind it, leaving implementation details to the following subsections.

Our solution to the problem described above is based on a fundamental observation: measuring precisely the RTT between all pairs of nodes (e.g., through pings) is extremely slow and costly, or even impossible due to topology dynamism. To



(a) SG-2 architecture.



(b) A superpeer topology in a 2-D virtual space.

Fig. 1. The figure on the left (Figure 1(a)) depicts the set of services composing the SG-2 architecture, while the figure on the right (Figure 1(b)) shows a superpeer topology in a bi-dimensional virtual space, where Euclidean distance corresponds to latency.

circumvent this problem, and allow nodes to estimate their latency without direct communication, the concept of *virtual coordinate service* has been developed [11]. The aim of this service is to associate every node with a synthetic coordinate in a virtual, n -dimensional space. The Euclidean distance between the coordinates of two nodes can be used to predict, with good accuracy, the RTT between them; in other words, it is sufficient for two nodes to learn about their coordinates to estimate their latency, without direct measurements.

Our problem may be redefined based on the concept of virtual coordinates. Nodes are represented by points in the virtual space; each of them is associated with an *influence zone*, described as a n -dimensional sphere of radius tol centered at the node. Our goal is to cover the virtual space with a small number of superpeers, in such a way that all nodes are either superpeers or are included in the influence zone of a superpeer. Figure 1(b) shows the topology resulting from the execution of SG-2 in a bi-dimensional virtual space.

Nodes communicate with each other using a *local broadcast service*, whose task is to efficiently disseminate messages to nodes included in the influence zone of the sender. This service is used by powerful nodes to advertise their availability to serve as superpeers, and by ordinary nodes to seek superpeers whose capacity has not been saturated yet.

The main component of SG-2 is the *superpeer management service*, which selects the superpeers and associates clients to them. The protocol is heavily inspired by the behavior of social insects [12], such as ants or bees, that have developed very sophisticated mechanisms for labor division. In summary, such mechanisms work as follows. In a totally decentralized fashion, specialized groups of individuals emerge, with each group aimed at performing some particular task. The task allocation process is dynamic and follows the community needs according to changes in the environment. The stimulus to perform some kind of task or to switch to another one can be given by many factors, but it is normally given by high concentrations of chemical signals, such as pheromones, that are released by other individuals and are spread in the environment. Each individual has its own response threshold to the stimulus and reacts accordingly.

The superpeer protocol mimics this general picture. Un-associated nodes diffuse a “request for superpeers” signal through local broadcasts; the signal concentration in the network may stochastically trigger a switch to the superpeer role in some nodes according to their response threshold, which is proportional to their capacity. On the other hand, powerful nodes covering the same area of the virtual space compete with each other to gain new clients, by signaling their availability through local broadcasts. Clients associate themselves to the most powerful superpeers, and superpeers with an empty client set switch back to the client role. The combination of these two trends (the creation of new superpeers to satisfy client requests and the removal of unnecessary superpeers) finds its equilibrium in a topology that approximates out target topology.

The last component to be described is the *peer sampling service*. The task of this layer is to provide each node with a view containing a random sample of nodes [13]. The motivation is twofold: first of all, the random sample is used by the local broadcast service to perform gossiping; second, the topology resulting from this layer can be described as a random graph composed of a single connected component among all nodes. This topology is extremely robust and present no central point of failure; it may be used to recover from catastrophic failures in the overlaying superpeer topology, for example due to a coordinated attack to the subset of superpeers.

A. Virtual Coordinate Service

In SG-2, the virtual coordinate service is provided by VIVALDI [11], which is a decentralized, scalable, and efficient protocol developed at MIT. Using VIVALDI, nodes may obtain good coordinates with few RTT probes directed to a small subset of nodes. More importantly, VIVALDI can exploit normal traffic produced by applications using it, without requiring further communication.

The *estimate* of the latency distance between v_i and v_j is denoted $est(v_i, v_j)$. Being estimates, these values may

differ from the actual latency. The pairwise error between the estimate and the actual latency can be computed as:

$$\frac{|lat(v_i, v_j) - est(v_i, v_j)|}{\min\{est(v_i, v_j), lat(v_i, v_j)\}}$$

In our experiments, the number of dimensions of the virtual space is 5; measuring the error between all pairs of nodes, we found a median error of only 0.14, and a maximum error of 3.5. Note that latency distances that are “under-estimated” may pose a problem: if the actual latency is over tol , but the estimated latency is smaller, a superpeer may accept a client out of the tolerated range. For this reason, the maximum error must be considered when selecting parameter tol .

B. Peer Sampling Service

The sampling service is provided by NEWSCAST [14], which has proven to be a valuable building block to implement several P2P protocols [15]. We provide here a brief description of the protocol and its characteristics.

Each NEWSCAST node maintains a view containing c node descriptors, each of them composed of a remote node identifier and a logical time-stamp. NEWSCAST is based on the gossip paradigm: periodically, each node (i) selects a random peer from its partial view; (ii) updates its local descriptor; and (iii) performs a *view exchange* with the selected peer, during which the two nodes send each other their views, merge them, and keep the c freshest descriptors.

This exchange mechanism has three effects: views are continuously shuffled, creating a topology that is close to a random graph with out-degree c ; the resulting topology is strongly connected (according to experimental results, choosing $c = 20$ is already sufficient for very stable and robust connectivity); and finally, the overlay topology is self-repairing, since crashed nodes cannot inject new descriptors any more, so their information quickly disappears from the system.

The peer sampling service is a key component both during the initialization phase (*bootstrap*) of the other layers, and during the normal functioning of the protocol, when it allows the discovery of “distant” or newly joined peers from the entire network. NEWSCAST is extremely inexpensive: messages are small, and the periodicity of view exchanges may be as low as one message per minute [14].

C. Local Broadcast Service

Unlike previous layers, based on existing protocols, the local broadcast service has adapted an existing protocol for the specific needs of SG-2 [16]. Each message m is associated with the sender identifier s_m and a radius parameter r_m . Message m is delivered to all those nodes that are within latency distance r_m from s_m , as estimated by VIVALDI. Hence, the name SPHERECAST.

The protocol may be described as follows. When a node either receives a message or wants to multicast a new one, it forwards it to its local *fan-out*. The fan-out of node v for message m is given by the subset of neighbors known to v that are potentially interested in the message, i.e. whose distance from s_m is not larger than r_m . SPHERECAST does not maintain

its own topology; instead, it relies on the underlying overlay network provided by the peer sampling service.

When a message is originated locally, or it is received for the first time, it is forwarded immediately to all nodes in the fan-out. If a message has been already received, a node may stochastically decide to drop it (i.e., not forwarding it). This is a standard approach used to avoid flooding the network. A strict deterministic approach such as dropping any multiple copy would not work correctly due to the nature of the underlying overlay. The actual clustering coefficient of the underlying topology and the continuous rewiring process may stop the message spreading. The stochastic approach solves this issue in a straightforward manner.

The probability of dropping a message is given by the following formula: $p = 1 - e^{-s/\vartheta}$, where s is the number of times the node has seen this message and ϑ is a response threshold parameter. In this way, when a packet is received multiple times by a peer, it has less and less probability to be forwarded again. From an implementation point of view, digests of received messages are stored in a per-node table, together with the number of times that specific message has been received. This table is managed with a LRU policy, to avoid unbounded growth.

D. Superpeer Management Service

This layer is the core component of SG-2. Nodes participate in this protocol either as superpeers or as clients; a client c may be either associated to a superpeer ($super(c) = s$), or actively seeking a superpeer in its tol range ($super(c) = \perp$). At the beginning, all nodes start as clients; to converge to the target topology defined in Section III, nodes may switch role at will, or change their client-superpeer relationship. The decision process is completely decentralized.

Each node v maintains the following local variables. *role* specifies the role currently adopted by v ; $role = SP$ if v is a superpeer, $role = CL$ otherwise. *clv* and *spv* are two views, respectively containing the clients and the superpeers known to v . They are composed of node descriptors combining an identifier w and a logical time-stamp ts_w ; the latter is used to purge obsolete identifiers, as in NEWSCAST. When v acts as a superpeer, *clv* is populated with the clients currently associated to v ; it is empty otherwise. The size of *clv* is limited by $cap(v)$. *spv* contains descriptors for the superpeers that are in $tol + \delta$ range; its size is not explicitly limited, but rather is bounded by the limited number of superpeers that can be found within $tol + \delta$ distance. When v acts as a client, one of the descriptors in *spv* may be the associated superpeer of v .

Two distinct kinds of messages are broadcasted using SPHERECAST: CL-BCAST and SP-BCAST. The former are sent while in client state and are characterized by a radius parameter r_m equal to tol , i.e. the maximum tolerated latency. The latter are used in superpeer state and their radius parameter is equal to $tol + \delta$; superpeers need a wider radius to get a chance to contact other superpeers; furthermore, nodes with overlapping influence zones can exchange clients if they find a better client allocation that reduces their latency.

At each node, two threads are executed, one active and one passive. The execution of active threads may be subdivided

TABLE I

BRIEF NODE BEHAVIOR DESCRIPTIONS ACCORDING TO EACH POSSIBLE PAIR OF NODE ROLE AND RECEIVED MESSAGE.

Node role	Received message structure	Behavior description
superpeer v	$\langle \text{SP-BCAST}, s, ts_s, cap(s) \rangle$	(1) insert or update (s, ts_s) in clv ; (2) if $cap(v) > cap(s)$, migrate s clients in v range to SP v until v capacity is exhausted; (3) if $ clv(v) =0$ et $est(v, s) \leq tol$, s joins v as a client
superpeer v	$\langle \text{CL-BCAST}, c, ts_c \rangle$	if the capacity of v has not ben exhausted, the client c joins v
client v	$\langle \text{SP-BCAST}, s, ts_s, cap(s) \rangle$	(1) insert or update (s, ts_s) in spv ; (2) if v has no SP, it asks s to be associated with him (theattpmt may fail if s has exhausted its capacity); (3) if v is already a client of superpeer s' and $cap(s) > cap(s')$, then it tries to migrate to the more powerful superpeer
\star client v	$\langle \text{CL-BCAST}, c, ts_c \rangle$	node v can change role according to the probability $p = \frac{s^2}{s^2 + \theta_v^2}$, where s is the signal magnitude and θ_v is the response threshold of node v . Check the text for the full explanation.

in periodic *cycles*: in each cycle, superpeers emit a SP-BCAST signal which is broadcast in the surrounding area, to notify nodes about their presence and its residual capacity. Clients, on the other hand, periodically emit CL-BCAST messages if and only if they are not associated to any superpeer. The shorter the cycle duration, the faster the system converge to the target topology; but clearly, the overhead grows proportionally. The passive threads react to incoming messages according to the message type and the current role as depicted in Table I. The SG-2 behavior is described in the third column, using a natural language like pseudo-code. Four distinct cases are possible, as the number of possible pairs of node role and message type received.

The fourth case, highlighted by the \star symbol in the table, is the cornerstone of our approach and needs a deeper explanation: this kind of message can trigger a role change from client to superpeer. The willingness of becoming a superpeer is a function of a node threshold parameter and the signal concentration perceived by a node in its influence area. The switching probability can be modeled by the following function:

$$P(\text{role}(v) = \text{CL} \rightarrow \text{role}(v) = \text{SP}) = \frac{s^2}{s^2 + \theta_v^2}$$

where s is the signal magnitude and θ_v is the response threshold of node v . This function is such that the probability of performing a switch is close to 1, if $s \gg \theta$, and it is close to 0 if $s \ll \theta$. If c_{max} is the maximum capacity, θ_v is initialized with a value which is $c_{max} - cap(v)$; in this way, nodes with higher capacity have a larger probability of becoming superpeers. The maximum capacity may be either known, or it can be easily computed by an *aggregation* protocol in a robust and decentralized fashion [15].

After the initialization, in order to make the topology more stable and avoid fluctuations, the response threshold is modified in such a way that time reinforces the peer role: the more time spent as a client, the less probable it is to change role. Once again, the inspiration for this approach comes from biology: it has been observed, for example, that the time spent by an individual insect on a particular task produces important changes in some brain areas. Due to these morphologic changes, the probability of a task change (e.g., from foraging to nursing) is a decreasing function of the time spent on the current task [12]. For this reason, θ_v is reinforced

as follows:

$$\theta_v(t) = \theta_v(t-1) + (\alpha \cdot (t - t'_v))$$

Where t is the current cycle and t'_v is the last cycle in which v became a superpeer; α is a parameter to limit or increase the time influence. The peer normal responsiveness is re-initialized based on its local capacity if its superpeer crashes or if it becomes a superpeer node.

The reaction to CL-BCAST messages is the only mechanism to allow a client to become a superpeer. A superpeer can switch back to the client role only when other higher capacity superpeers have drained its client set. The θ adaptation process is only active when a node is in the client state.

V. EXPERIMENTAL RESULTS

We performed a large number of experiments based on simulation to validate the effectiveness of our approach. The goal of our experiments was twofold: first of all, we measured the speed of convergence in a stable overlay, in the absence of failures; second, we measured the robustness of our approach in a dynamic environment, where a fixed percentage of nodes are substituted with fresh ones periodically. Finally, communication overhead has been measured. The experiments have been performed using PeerSim [17].

In our experiments, network size is fixed at 1000 and 2000 nodes. Several kinds of networks have been considered, but here, the focus is on a *gaming-oriented* scenario [18], [19]; in addition, we also present some preliminary results regarding a *P2P IP-telephony* scenario. The network size for the second scenario is small compared, for example, with the millions of Skype users, but we are limited by the simulations constraints and by the usage of real world latency data.

The approach to build our virtual coordinates distribution (see Section IV-A) follows a common off-line procedure for both scenarios: (i) we build a matrix holding for each pair of nodes v, w the latency distance $lat(v, w)$ between them; (ii) the Vivaldi [11] algorithm is run on the latency matrix obtaining the corresponding function $est(v, w)$ for each matrix location.

The way we have obtained the latency matrix varies according to the actual scenario adopted. In the gaming scenario, supported by the results in [18], [20], we adopted a normal distribution with average value $\mu = 250 \text{ ms}$ and variance $\sigma = 0.1$. In the IP-Telephony scenario instead, we used the

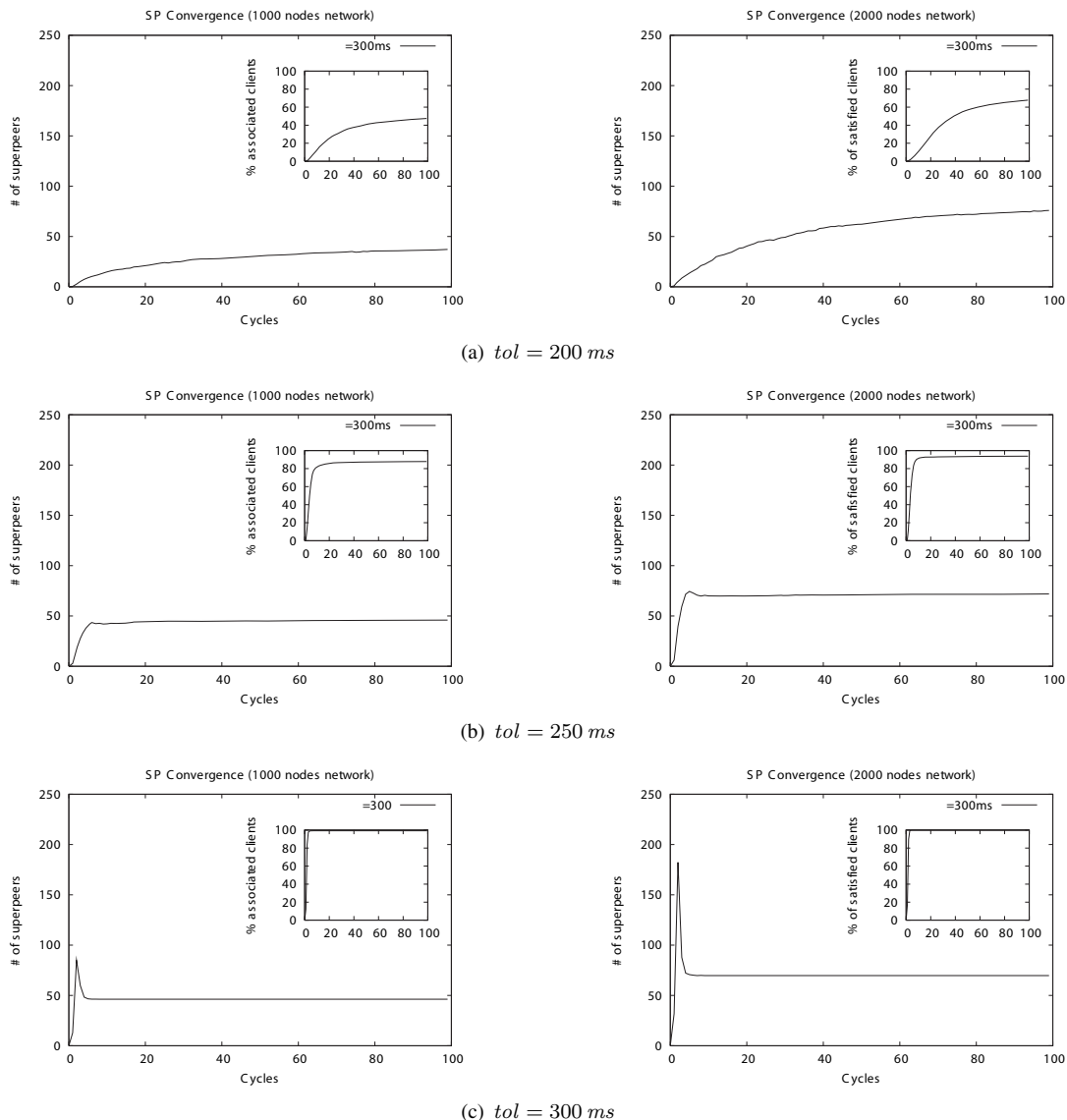


Fig. 2. Convergence time. Three tol values are considered: $200\ ms$ (a), $250\ ms$ (b), $300\ ms$ (c). The main figures show the number of active superpeers at each cycle, while the small sub-figures show the number of clients that are in tol range.

Meridian data-set [21] as a latency matrix. This data set is a sample of real world Internet latencies among 2500 nodes spread all over the world. We decided to build the coordinates off-line in order to speed up considerably the simulations. The simulation is organized in synchronous *cycles*, during which each node has the possibility to initiate a gossip exchange; note, however, that in reality nodes do not need to be synchronized. This fact allows the application designers to set the cycle length to any wall clock time interval they need. In both scenarios, we have experimented with distinct δ values in the range $[200 : 400]\ ms$, corresponding to typical round-trip time that can be accepted for superpeer communication. We decided to include only the results for $\delta = 300\ ms$ due to its small average advantage shown in our simulations. The capacity function $cap()$, i.e. the maximum number of clients that can be served, is generated through an uniform distribution in the range $[1 : 500]$. All the results are averaged over 10 experiments.

A. Gaming scenario discussion

We assume that each SP manages a replica of the global game state. The actual synchronization mechanism is an application dependent detail and we do not address it in this work. In the corresponding virtual space, we have considered tol values of $200\ ms$, $250\ ms$ and $300\ ms$, which are typical of strategy and role-playing games.

Overlay convergence

Figure 2 illustrates the behavior of the protocol over time. All the figures in the left column are obtained in networks whose size is 1000 nodes, while the figures in the right column are relative to networks with size equal to 2000 nodes. The content of each sub-figure is divided in two parts; in the main plot, the number of superpeer active at each cycle is shown; in the small frame inside the main plot, the percentage of clients that are already associated is shown. In these experiments, the network is static; no nodes are removed or added.

Figure 2(a) depicts a rather bad situation: in both network sizes, the convergence is extremely slow, and the number

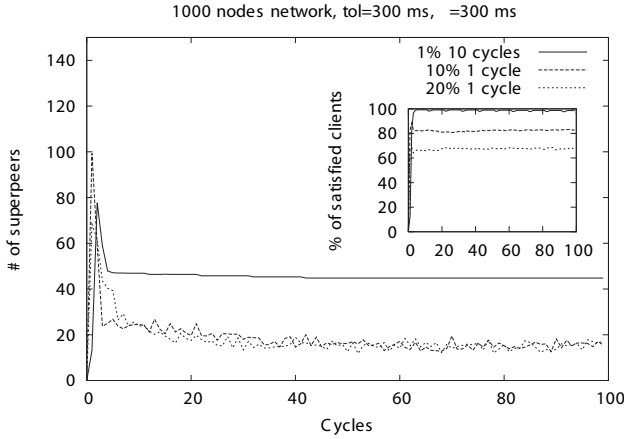


Fig. 3. Experiments with churn. Network size is 1000; periodically, 1% (every 10 cycles), 10% or 20% (every cycle) of the nodes are substituted with new ones.

of nodes that are satisfied is low. This bad performance is motivated by the characteristics of the latency distributions [18], [19] and the tolerance value selected; most of the node pairs have a higher latency than 200 ms , and thus SG-2 cannot help much. Figure 2(b) shows a much better situation: a large percentage of clients (between 94% and 100% depending on size) have been associated after only few cycles (10-20). The number of superpeers is also very small, after an initial peak due to a large number of clients reacting to the signal. Almost every client can reach the required latency because 250 ms is the average pairwise latency in our game-like coordinates distribution. However, some nodes lie outside the 250 ms border and it is challenging for SG-2 to accommodate those nodes. The node density plays an important role for SG-2. In fact, the bigger population produces a higher stimulus concentration that produces a faster reaction; the network can be organized in a (almost stable) latency-aware fashion in just 20 cycles.

Figure 2(c) shows the performance for 300 ms tol : a response time that is perfectly acceptable in a strategic/role playing game scenario. We obtain 99.9% of in range clients with about 50 superpeers and 100% with about 63 superpeers respectively in the small and bigger network; in both cases the time required is less than 10 cycles.

Churn

Figure 3 is aimed at illustrating the robustness of our protocol. The size of the network is fixed at 1000 nodes. Its composition, however, is dynamic: periodically, a fixed set of the peer population crashes and it is substituted with new peers. Any node in the network can be affected by substitution, regardless of its role. Unlike the real world, where a superpeer is supposed to be more reliable, our choice is stricter and more “catastrophic”. We consider three distinct set sizes: 1%, 10%, 20% of the network size. The first set is substituted every 10 cycles, while the others are substituted at each cycle. The last two churning sets are very large and un-realistic, but they are ideal candidates to show the SG-2 performance in a disastrous scenario. The first set instead, is realistic and taken from real world churn measurements in a DHT environment [22]. To

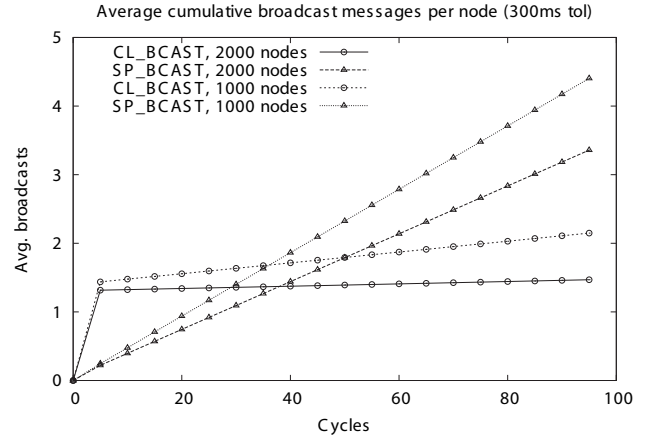


Fig. 4. Broadcast message traffic; the traffic is made by CL-BCAST and SP-BCAST message types. Both network sizes results are shown.

conform fairly to the results in [22], we consider the temporal size of the cycle equal to 1 minute.

Using the “catastrophic” values, Figure 3 shows that the number of superpeers oscillates over time, as expected, and that up to 80% and 70% of the clients are associated to superpeers. The nodes that are not associated are those that have been recently created and are trying to find a position in the topology.

The realistic value instead, does not affect the latency-aware topology. In fact, after cycle 42, the superpeer population is fully stable and its size is comparable to Figure 2(c).

Message overhead and message loss

Finally, we discuss the broadcast message overhead and the performance in case of message loss.

In Figure 4, we have measured the average number of diffused broadcast messages, distinguishing among CL-BCAST and SP-BCAST. The system *tol* parameter is set to 300 ms as in Figure 2(c). Since the CL-BCAST message type is broadcast only in case of lack of satisfaction, only a small number of them are generated: on average, less than 2 messages if the SG-2 can fully optimize the overlay (as in the bigger network shown in Figure 2(c)). Superpeers, on the other hand, continuously send one broadcast message per cycle. The bigger network has an advantage since the percentage of superpeers needed to optimize the network is much less than that in the smaller overlay, therefore the average number of SP-BCAST per node is about 30% lower. This fact emphasizes again how SG-2 is dependent on the virtual space density.

Finally, we run our superpeer architecture over an unreliable transport layer. We consider a 3% loss probability on any message exchange; SPHERECAST is the actual protocol layer affected by the message loss. Figure 5 compares the impact of the message loss versus the “perfect” (e.g., without loss) transport scenario on a convergence performance basis. The comparison regards both network sizes. As expected, SG-2 takes longer to converge to the optimized overlay. In the bigger network, it takes more than twice the previous time (about 22 versus 8 cycles). The amount of satisfied clients, depicted in the small sub-figures, is the same at the expense of a longer time to wait; on the other hand, the number of superpeer

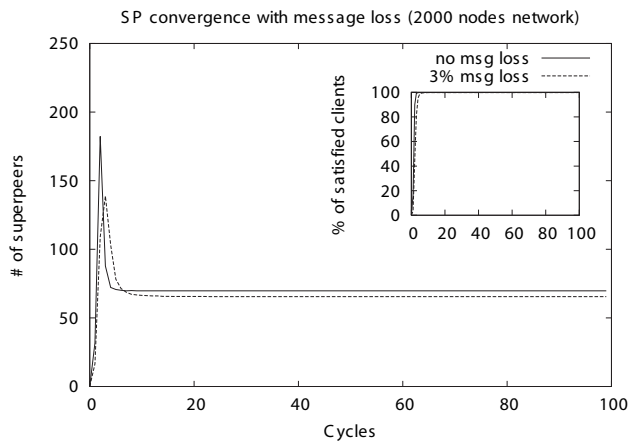


Fig. 5. Experiments with message loss. Each message has a 3% probability to be lost by the transport layer. The network size is 2000 nodes and the tol value is 300 ms. The main figure shows the number of active superpeers over time, while the small sub-figures show the % amount of clients that are in the tol range.

required is slightly lower than in the case of the absence of message loss. This fact is due to the longer time taken to converge: it makes less probable to fall in a local minimum (see Section III) and allows a more coordinated emergence of superpeer nodes.

B. IP-Telephony scenario discussion

We consider a particular superpeer-based IP-Telephony (IPT) system similar to Skype [4], but designed over open standards such as the SIP protocol [9], [10] and we do not address the interaction with existing PSTN networks. The superpeer design makes easier the location of contacts over the overlay and provides a firewall piercing capability. A caller locates the desired contact through the facilities provided by its direct SP, then SP act as a router delivering the compressed audio speech between its client and the contacted node. In many cases, the communication would be faster if both parties could communicate directly among them; however this is not always a viable solution especially in a P2P overlay (e.g., due to NAT routing or firewall issues). This routing behavior is similar to Skype [7].

In this kind of applications the perceived QoS by the users is mainly influenced by the 1-way latency between them rather than the RTT. In general, the 1-way latency limit, under which the user QoS perception is not affected, is 150 ms (see [23]). This is why we considered tol values (RTT) of 300 and 400 ms for the superpeer overlay construction.

Figure 6 shows the emergence of SP nodes in the IP-Telephony like scenario; the results for both network sizes are shown.

In the smaller network, about 30 superpeers are required to ensure a QoS of at most 200ms 1-way latency for all the clients, while ensuring a limit of 150 ms requires about 55 super nodes. However, non all clients can acquire this QoS due to the bandwidth constraint and the latency distribution of the underlying network, but a 95% of satisfied clients can still be considered a good compromise. In the bigger network the

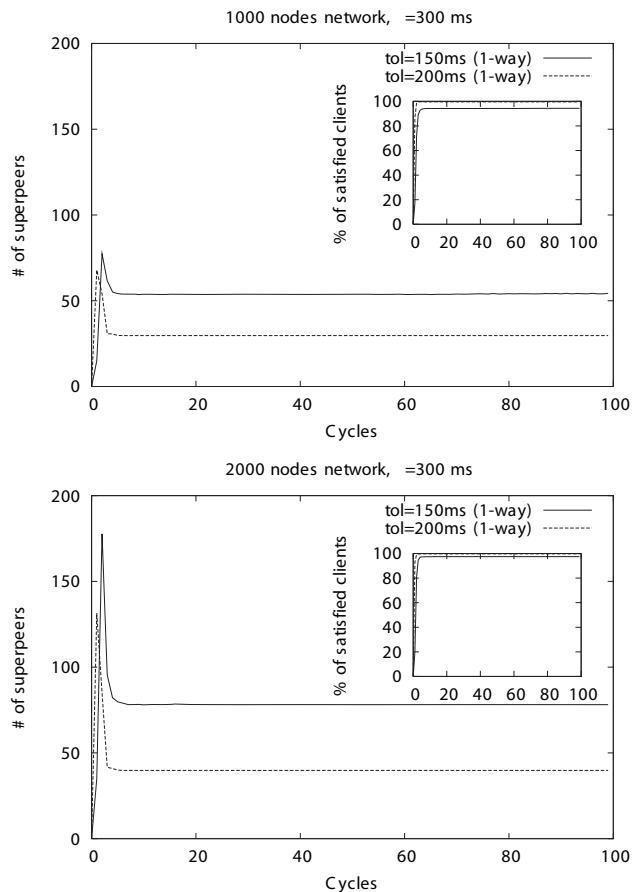


Fig. 6. Convergence to the SP topology in the IP-Telephony latency model. The results for both network sizes are shown; the tol value are 150 and 200 ms 1-way latency. The main figures show the number of active superpeers over time, while the small sub-figures show the % of clients that are in the tol range.

situation is similar; of course more super nodes are required to satisfy the higher concentration of participants: about 41 and 80 SP respectively in the 200 and 150 ms case. In both networks, the convergence to a stable SP overlay is achieved in less than 10 gossip rounds.

In order to quantify the effectiveness of our IPT overlay we made a particular test in which we simulate a call among every pair of clients in the system. For each distinct pair of ordinary nodes (a, b) , the caller node a checks which is the SP_i in its range that minimizes the value of $lat(a, b)$, where $lat(a, b) = lat(a, SP_i) + lat(SP_i, b)$. Notice that a node holds in cache the IDs of any SP in the tol range. In other words, we try to find the “straightest way” to reach node b according to the available SP links.

The results achieved are summarized in Table II. The minimum, maximum and average 1-way latency among all call pairs is reported according to the network size and tol value adopted. The average latency is close to the optimal tolerable value (i.e., 150 ms) and becomes slightly higher when the tol parameter is 200 ms. The maximum values obtained are interesting: they are far from being optimal, but also far from being unusable. Of course, this is a simple scenario that does not consider many additional details, (e.g, such as the latency

TABLE II

IP-TELEPHONY SCENARIO: MANY-2-MANY CALL EXPERIMENT. ALL LATENCIES AND TOLERANCES ARE CONSIDERED 1-WAY AND ARE EXPRESSED IN *ms*.

Network size	tol	Min lat.	Max lat.	Avg. lat.
1000	150	30	339.64	158.2
1000	200	41.7	355.3	163.6
2000	150	23.95	327.19	152
2000	200	30.4	344.1	158.43

introduced by the SP while routing the packets or the audio buffering delay), but we consider these results encouraging.

VI. RELATED WORK

The superpeer approach to organize a P2P overlay is a trade-off solution that merges the client-server model relative simplicity and the P2P autonomy and resilience to crashes. The need for a superpeer network is mainly motivated by the fact to overcome the heterogeneity of peers deployed on the Internet.

Yang and Garcia Molina [24] proposed some design guidelines. A mechanism to split node clusters is proposed and evaluated analytically, but no experimental results are presented.

Superpeer solutions proved to be effective solutions in the real world: Kazaa / Fasttrack [2] and Skype [7] are two outstanding examples. However, their actual protocols are not publicly available and they cannot be compared with any other solution or idea. At the time of writing, only a few works [4], [25] describe some low-level networking details.

The SG-2 protocol can be considered as a natural evolution of SG-1 [26]; the two solutions, however, cannot be directly compared from a performance point of view because of their different goals. SG-1 focuses on optimizing the available bandwidth in the system, while SG-2 introduces the notion of latency between peer pairs and poses a QoS limit on it. The definition of the target topology is straightforward in SG-1 (e.g., the minimum number of superpeers to accommodate all the peers according to the superpeer capacities), while it is a NP-problem in the SG-2 case. From the architectural point of view, they both rely on the existence of an underlying random overlay, on top of which the superpeer overlay is generated. The superpeer election process in SG-2 is strongly bio-inspired and much more randomized than approach used in SG-1.

In [27], the authors propose a socio-economic inspiration based on Shelling's model to create a variation of the superpeer topology. Such variation allows ordinary peers to be connected with each other and to be clients of more than one superpeer at the same time. This topology focuses on efficient search. As in our case, the superpeers are connected to each other to form a network of hubs and both solutions are suited for unstructured networks. However, they do not address the problem of the superpeer election.

In [28], the authors suggest an interesting use of the aggregation protocol [15] in order to build a particular topology (the *gradient* topology). The discovery of superpeers nodes in such topology is a straightforward task. Their algorithm relies on the existence of an *utility* function that captures the peer application specific constraints. An utility value above

a certain threshold makes a node eligible for superpeer role. The interesting aspect of this work is that it is a very general superpeer framework in which SG-2 may fit (given a suitable utility function).

The basic problem of finding the best peer, having the required characteristics, to accomplish some task (e.i., the superpeer task) is addressed in a more general form in [29]. The problem is referred as "optimal peer selection" in P2P downloading and streaming scenarios. The authors use an economics inspired method to solve the optimization problem; the developed methodologies are general and applicable to a variety of P2P resource economy problems. The proposed solution is analytically strong, but no experimental results are shown especially regarding a large and dynamic scenario as the one the authors are addressing.

Our implementation is based on VIVALDI (see section IV-A), but it is not tied to any particular virtual coordinate service. Other architectures can be adopted, such as IDMaps [30] and GNP [31], or PIC [32] and PCoord [33]. The first two rely on deployment of infrastructures nodes, while the others provide latency estimates gathered only between end-hosts, as VIVALDI does. We opted for VIVALDI because of its fully distributed nature and simple implementation.

In less strict latency context, such as file sharing, hop count is usually preferred in contrast to actual latency to provide distance estimation. Pastry [34], [35], for example, uses a hop distance metric to optimize its response time.

VII. CONCLUSIONS

This paper presented SG-2, a fully decentralized, self-organizing general protocol for the construction of proximity-aware, superpeer-based overlay topologies. The protocol produces an overlay in which almost all nodes (99.5%) are in range with a *tol* latency of 300 *ms*. The number of generated superpeers is small with respect to the network size (only 3-5%) in both test scenarios. The protocol shows also an acceptable robustness to churn.

We conclude noting that the results presented in this paper are only a first step toward the implementation of fully decentralized P2P latency-aware applications; several other problems have to be solved, including security, state replication, state distribution, etc. Our current efforts are mainly focused on security issues.

ACKNOWLEDGMENT

This work was partially supported by the FET unit of the European Commission through projects BISON (IST-38923), DELIS (IST-01907), CASCADAS (IST-27807) and BIONETS (IST-27748).

REFERENCES

- [1] "Gnutella web site," <http://gnutella.wego.com>.
- [2] "Fasttrack web site," <http://www.fasttrack.nu>.
- [3] P. Bettner and M. Terrano, "1500 archers on a 28.8: Network programming in Age of Empires and beyond," in *Proc. GDC'01*, Mar. 2001.
- [4] S. Baset and H. Schulzrinne, "An analysis of the Skype peer-to-peer internet telephony protocol," Columbia University, Department of Computer Science, New York, NY, Tech. Rep. CUCS-039-04, Sept. 2004.
- [5] "PeerCast p2p radio," <http://www.peerCast.org>.

- [6] L. Gautier and C. Diot, "MiMaze, a multiuser game on the internet," INRIA, France, Tech. Rep. RR-3248, Sept. 1997.
- [7] "Skype: Free internet telephony that just works," <http://www.skype.com>.
- [8] [Online]. Available: <http://www.cs.columbia.edu/sip/>
- [9] D. Bryan and C. Jennings, "A p2p approach to sip registration and resource location," July 2005. [Online]. Available: <http://www.p2psip.org/drafts/draft-bryan-sipping-p2p-01.html>
- [10] [Online]. Available: <http://www.p2psip.org/>
- [11] F. Dabek, R. Cox, F. Kaashoek, and R. Morris, "Vivaldi: A decentralized network coordinate system," in *Proc. SIGCOMM '04*, Aug. 2004.
- [12] E. Bonabeau, M. Dorigo, and G. Theraulaz, *Swarm Intelligence: From Natural to Artificial Systems*. New York: Oxford University Press, Inc., 1999.
- [13] M. Jelasity, R. Guerraoui, A.-M. Kermarrec, and M. van Steen, "The peer sampling service: Experimental evaluation of unstructured gossip-based implementations," in *Proc. 5th Int. Middleware Conf.*, Oct. 2004.
- [14] M. Jelasity, W. Kowalczyk, and M. van Steen, "Newscast computing," Vrije Universiteit Amsterdam, Department of Computer Science, Amsterdam, The Netherlands, Tech. Rep. IR-CS-006, Nov. 2003.
- [15] M. Jelasity, A. Montresor, and O. Babaoglu, "Gossip-based aggregation in large dynamic networks," *ACM Trans. Comput. Syst.*, vol. 23, no. 1, pp. 219–252, 2005.
- [16] P. Eugster, R. Guerraoui, S. B. Handurukande, A.-M. Kermarrec, and L. Massoulié, "Lightweight probabilistic broadcast," *ACM Trans. Comput. Syst.*, vol. 21, no. 4, pp. 341–374, 2003.
- [17] "PeerSim peer-to-peer simulator," <http://peersim.sf.net>.
- [18] S. Zaniolas and R. Sakellariou, "Towards a monitoring framework for worldwide grid information services," in *Euro-Par*, 2004, pp. 417–422.
- [19] N. Sheldon, E. Girard, S. Borg, M. Claypool, and E. Agu, "The effect of latency on user performance in Warcraft 3," in *Proc. 2nd Workshop on Network and System Support for Games*. New York: ACM Press, 2003, pp. 3–14.
- [20] M. Allman, "A web server's view of the transport layer," *Comput. Commun. Rev.*, vol. 30, no. 5, pp. 10–20, 2000.
- [21] "Meridian internet data set." [Online]. Available: <http://www.cs.cornell.edu/People/egs/meridian/data.php>
- [22] R. Mahajan, M. Castro, and A. Rowstron, "Controlling the cost of reliability in peer-to-peer overlays," in *Proc. IPTPS'03*, Feb. 2003.
- [23] A. Percy, "Understanding latency in ip telephony." [Online]. Available: www.aitel.hist.no/fag/ipt/lek02/iptel_latency_brooktrout.pdf
- [24] B. Yang and H. Garcia-Molina, "Designing a super-peer network," in *Proc. IEEE International Conference on Data Engineering (ICDE'03)*, 2003.
- [25] N. Leibowitz, M. Ripeanu, and A. Wierzbicki, "Deconstructing the kazaa network," in *Proc. 3rd IEEE Workshop on Internet Applications (WIAPP'03)*, June 2003.
- [26] A. Montresor, "A robust protocol for building superpeer overlay topologies," in *Proc. of the 4th Int. Conf. on Peer-to-Peer Computing*. Zurich, Switzerland: IEEE, Aug. 2004.
- [27] A. Singh and M. Haahr, "Creating an adaptive network of hubs using Schelling's model," *Commun. ACM*, vol. 49, no. 3, pp. 69–73, 2006.
- [28] J. S. and Jim Dowling, R. Cunningham, and R. Meier, "Using aggregation for adaptive super-peer discovery on the gradient topology," in *SelfMan*, ser. Lecture Notes in Computer Science, A. K. and Jean-Philippe Martin-Flatin, Ed., vol. 3996. Springer, 2006, pp. 73–86.
- [29] M. Adler, R. Kumar, K. W. Ross, D. Rubenstein, T. Suel, and D. D. Yao, "Optimal peer selection for p2p downloading and streaming," in *Proc. IEEE Infocom'05*, March 2005.
- [30] P. Francis, S. Jamin, C. Jin, Y. Jin, V. Paxson, D. Raz, Y. Shavitt, and L. Zhang, "IDMaps: a global internet host distance estimation service," in *Proc. IEEE Infocom'99*.
- [31] T. Ng and H. Zhang, "Predicting internet network distance with coordinates-based approaches," in *Proc. IEEE Infocom'02*.
- [32] M. Costa, M. Castro, A. Rowstron, and P. Key, "Pic: practical Internet coordinates for distance estimation," in *Proc. ICDCS'04*.
- [33] L. Lehman and S. Lerman, "Pcoord: network position estimation using peer-to-peer measurements," in *Proc. 3rd IEEE International Symposium on Network Computing and Applications (NCA'04)*.
- [34] A. Rowstron and P. Druschel, "Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems," in *Proc. IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)*, Nov. 2001, pp. 329–350.
- [35] M. Castro, P. Druschel, Y. C. Hu, and A. Rowstron, "Exploiting network proximity in distributed hash tables," in *Proc. International Workshop on Future Directions in Distributed Computing (FuDiCo)*, O. Babaoglu, K. Birman, and K. Marzullo, Eds., June 2002, pp. 52–55.

PLACE
PHOTO
HERE

Gian Paolo Jesi is a Ph.D. student in computer science at the University of Bologna. His research interests are focused on large-scale P2P systems, emergent behavior, epidemic protocols and application of bio-inspired approaches to distributed systems problems. He received his MSc in computer science from the University of Bologna in 2002; after a short work experience in an Italian IT company, he had the opportunity to move to the University of Bologna joining the BISON project.

PLACE
PHOTO
HERE

Alberto Montresor is Associate Professor of Computer Science at the University of Trento, Italy. He received his Ph.D. in 2000 from the University of Bologna, Italy, where he designed Jgroup, a partition-aware group communication system. He served as Research Associate in Bologna until 2005, when he moved to Trento. He is author of more than 40 papers in international conferences and journals, and he has been active in several European projects in the field of distributed computing and complex adaptive systems.

PLACE
PHOTO
HERE

Ozalp Babaoglu is Professor of Computer Science at the University of Bologna, Italy. He received a Ph.D. in 1981 from the University of California at Berkeley where he was a principal designer of BSD Unix. He is the recipient of 1982 Sakrison Memorial Award, 1989 UNIX International Recognition Award and 1993 USENIX Association Lifetime Achievement Award for his contributions to the UNIX system community and to Open Industry Standards. Before moving to Bologna in 1988, Babaoglu was an Associate Professor in the Department of Computer Science at Cornell University. He is active in several European research projects in distributed computing and complex adaptive systems. Babaoglu is an ACM Fellow and serves on the editorial boards for ACM Transactions on Computer Systems, ACM Transactions on Autonomous and Adaptive Systems and Springer-Verlag Distributed Computing.