

Decentralized Schemes for Size Estimation in Large and Dynamic Groups

Dionysios Kostoulas, Dimitrios Psaltoulis, Indranil Gupta

Dept. of Computer Science

University of Illinois, Urbana-Champaign

{kostoula, psaltoul, indy}@cs.uiuc.edu

Ken Birman, Al Demers

Dept. of Computer Science

Cornell University

{ken, ademers}@cs.cornell.edu

Abstract— Large-scale and dynamically changing distributed systems such as the Grid, peer-to-peer overlays, etc., need to collect several kinds of global statistics in a decentralized manner. In this paper, we tackle a specific statistic collection problem called *Group Size Estimation*, for estimating the number of non-faulty processes present in the global group at any given point of time. We present two new decentralized algorithms for estimation in dynamic groups (i.e., where processes are constantly joining and crashing), analyze the algorithms, and experimentally evaluate them using real-life traces. One scheme is active: it spreads a gossip into the overlay first, and then samples the receipt times of this gossip at different processes. The second scheme is passive: it measures the density of processes when their identifiers are hashed into a real interval. Both schemes have low latency, scalable per-process overheads, and provide high levels of probabilistic accuracy for the estimate. They are implemented as part of a size estimation utility called PeerCounter that can be incorporated modularly into standard peer-to-peer overlays. We present experimental results from both the simulations and the PeerCounter utility, running on a cluster of 33 Linux servers. The studies measure and compare microbenchmarks for these approaches in static groups, as well as in dynamic groups with significant turnaround (arrival and departure) of processes – the latter experiments are driven by real-life traces.

I. INTRODUCTION

Distributed systems such as peer-to-peer overlays, sensor networks, Grid application overlays, etc., tend to be large-scale since they contain several thousands of processes. More importantly, however, they are also *dynamic*. Dynamism means that there is continuous arrival and departure activity through processes joining, crashing and voluntarily departing. At the same time, distributed applications often require an estimate of the number of non-faulty processes currently present in the group. We call this as the problem of *Group Size Estimation*. The problem is challenging not only due to the dynamism and scale involved, but also because the above mentioned peer-to-peer overlays typically require each process to maintain only *partial* information concerning the group membership.

Accurate estimates of the number of processes currently present in the overlay, i.e., the group size, are absolutely essential in order to enable several notions of run-time adaptivity in peer-to-peer applications. For example, in a Grid application, continuous estimates of the group size can be used to dynamically partition a distributed problem among participating clients. In peer-to-peer distributed hash tables such as Pastry [14] and

Chord [15], the estimated value is required to set the routing table size and to set timeouts for queries – the number of virtual hops for routing a query depends on the group size. Other overlays, e.g., Kelips [9], use the notion of process subgroups whose size depends on the number of processes present in the system.

The Group Size Estimation problem is representative of a large class of problems for collecting *statistics* about a large-scale distributed system, in a decentralized manner. Other problems in this class include aggregation, e.g., [3]. However, the specificity of the Group Size Estimation problem lends itself to solutions with the potential for much greater accuracy and scalability than that obtained by the straightforward application of aggregation algorithms.

The Group Size Estimation problem has two flavors – one-shot and continuous. We formally define the one-shot problem next – the continuous version is similar.

Group Size Estimation Problem: *Give a protocol that when initiated by one process, estimates the number of non-faulty processes in the overlay graph component containing the initiating process.*

Impossibility of Group Size Estimation in a Dynamic Group: This problem is impossible to solve accurately in a dynamic group. Notice that a group size estimation protocol will take non-zero time to run, and for any process p that is not the initiator, there will be a non-zero delay between p 's last message in the estimation protocol, and the initiator finalizing the estimate. In a run where process p fails during this interval, the estimate will be incorrect.

This motivates the design of algorithms for *approximate* estimation. One simple algorithm could be the following. The initiator process sends a multicast to the overlay through a dynamic spanning tree, and waits to receive back replies. The disadvantage of this scheme is that it is non-fault-tolerant – if a process fails after receiving the multicast and before replying to its parent in the spanning tree, the estimate will exclude all descendants of that process in the tree.

Contributions of the Paper: This paper proposes practical, decentralized, efficient, and fault-tolerant estimation algorithms with a probabilistic output. Specifically, we study two algorithms, both based on variants of sampling. The first algorithm is active: called the *Hops Sampling Algorithm*, it initiates a gossip into the

group, and samples, based on hop distance, the times at which the gossip is received by different processes. The second approach is passive: called the *Interval Density Approach*, it samples the density of processes that lie in a given real interval, when the process identifiers are hashed.

Both our algorithms can be run either in a one-shot manner or on a continuous basis (the latter variant is preferable for long-running applications such as peer-to-peer overlays). When used in a one-shot manner, they have running times that grow logarithmically with group size, and impose a small sublinear overhead on each process in order to achieve an estimate that is accurate with high probability (w.h.p.). In the continuous version, the time taken for a given process arrival/failure/departure to affect the group size estimate, grows logarithmically with the group size.

We present experimental evaluation of the two algorithms. The first set of experimental results is from a simulation, and includes micro-benchmarks and trace-based experiments. The second set of experiments uses *PeerCounter*, our open-source implementation of the size estimation algorithms. Using *PeerCounter* allows our simulations to be run on a cluster of up to 33 Linux servers, each with multiple processes. *PeerCounter* is open-source and can be easily incorporated into a variety of peer-to-peer substrates, e.g., Pastry, Chord, Kelips [9, 14, 15], etc., thus providing them with the capability to estimate the number of non-faulty processes present in the group.

System Model: We assume a group of processes with unique identifiers, communicating through an asynchronous network, and connected in an overlay. Messages can be sent from a process u to any other process v whose identifier is known by u at its local membership list – these membership entries are nothing but the links in the peer-to-peer overlay. Membership entries are assumed to be maintained by a complementary membership protocol already running within the overlay. The membership protocol needs to be only weakly consistent, i.e., detect failures eventually. It is not required to provide strong guarantees such as virtual synchrony. We detail other requirements from the membership protocol where needed.

The number of non-faulty processes is N , and is an unknown quantity. Our protocols operate in *rounds* at each process, with the round duration fixed across the group. Processes are not required to have synchronized

clocks. However, since the duration of a round is $O(\text{seconds})$, processes can be assumed to have negligible clock drifts. Processes can undergo crash-stop failures – crash-recovery failures can be supported by having a process rejoin the group with an identifier that is unique and unused previously.

The rest of this paper is organized as follows. Section II discusses related work. Sections III and IV present the Hops Sampling algorithm and the Interval Density approach respectively. Sections V and VI give experimental results from our implementations. We summarize in Section VII.

II. RELATED WORK

Several efforts have been made to provide solutions to the group size estimation problem. Unfortunately, most of these approaches either do not apply to dynamic groups, or make unreasonable assumptions. Ref. [3] presents several mechanisms for aggregation and group size estimation. One algorithm has an estimation message sent along a random walk within the group. When the estimation message hits a process for the second time, the protocol stops and the number of hops traversed so far can be used to estimate the group size – according to the birthday paradox, it takes $\Theta(\sqrt{N})$ hops for a random walk message to encounter a process a second time. However, the latency of this algorithm is $\Theta(\sqrt{N})$, and this is impractical in groups that have processes joining and leaving at high rates. Our schemes have a latency that increases only logarithmically with N .

Our Interval Density algorithm bears resemblance to the approach in Ref. [1], but there are significant differences. Ref [1] uses a hash function to assign process IDs in the real interval $[0,1]$ and decomposes the $[0,1]$ interval into a hierarchy of subintervals, each holding identifiers from processes that run on the same region. A recursive algorithm is used, starting from the peer’s home region and continuing upwards until an estimate for the $[0,1]$ interval is achieved. In comparison, our Interval Density approach does not depend on the way process identifiers are assigned in the group. Further, the algorithm of [1] could suffer from message implosion at the initiator process, while the Interval Density algorithm imposes no extra message complexity since it is a passive approach.

Kermarrec et al also discuss practical estimation schemes in [11], but at the time of writing this paper, they have not been evaluated. A ring-based algorithm for group size estimation is discussed in [12]; however, the estimated group size can have a large error – it can be between $N/2$ and N^2 , where N is the actual group size. The required logical ring may not be present in all overlays, e.g., Grid overlays. Our schemes in this paper do not require such a fixed minimal overlay, and can be run on a large range of overlays.

Epidemic algorithms were discussed by Demers et al in [6], Eugster et al in [7], Karp et al in [10], and originate from the Mathematical study of epidemics [2]. Birman et al [5] use gossip to design a probabilistically reliable multicast protocol.

III. ACTIVE APPROACH - THE HOPS SAMPLING ALGORITHM

Figure 1 shows pseudocode for the canonical Hops Sampling algorithm. The algorithm is proactive, operating in *protocol periods* (also called *rounds*). The duration of a *round* is fixed and a round starts at the same time at all processes. The estimation is started by one process called the *initiator*, but the rest of the protocol operates in a decentralized manner throughout the group.

The initiator sends an *initiating message* at the start of the protocol. Once a process has received the initiating message, at the beginning of each subsequent round, it selects `gossipTo` other processes as targets and sends them *gossip* messages containing the initiating message. A process stops gossiping once either `gossipFor` rounds have expired since receipt of the initiating message, or `gossipUntil` gossip messages have been received by the process. All `gossip*` parameters mentioned above are a priori fixed integer constants. Besides the local membership list, each process p also maintains a list, `fromList`, of some other processes that *it* knows to have already been infected – these are simply the processes that have sent gossip messages to process p . The selection of gossip targets by a process is done *uniformly at random* from the membership list, but by excluding the elements that appear in `fromList`.

Since the membership protocol is weakly consistent, the membership list may be out of date. In order to maintain the `gossipTo` parameter, gossip messages need to be acknowledged, and a gossiping node retries each gossip message with a different target until it is acknowledged.

Hops Sampling Protocol (gossipTo, gossipFor, gossipUntil, gossipResult, gossipSample):

/ Single initiator, Single instance. Multiple runs will require extra space and processing for a unique run number. */*

Gossip message format: (hopNumber, initiator, from);

Group Membership List: G;

fromList = NULL; / list of processes from which I have received gossip */*

At the Initiator process p :**Initially:**

gossip_on = TRUE; // am I actively gossiping?

past_gossip_rnds = 0; // number of past local gossip rounds

myHopCount = 0;

num_rcvdgossips = 0; // number of gossips received so far

initiator = p ;

Once every protocol period:

if (*gossip_on == TRUE*)

send message (l, p, p) to gossipTo processes selected uniformly at random from set $G - \text{fromList}$;

past_gossip_rnds++;

if (*past_gossip_rnds > gossipFor*)

gossip_on = FALSE;

On receipt of a gossip message M : */* received messages processed serially */*

fromList = fromList + { $M.\text{from}$ };

num_rcvdgossips++;

if (*num_rcvdgossips > gossipUntil*)

gossip_on = FALSE;

gossipResult protocol periods after starting:

sample a random subset of gossipSample processes for their myHopCount values;

for each non-zero value received do: HopReportArray[myHopCount]++;

output the mean of received non-zero values as:

(weighted sum of different values) / gossipSample;

At a non-initiator process p :**Initially:**

gossip_on = FALSE; // am I actively gossiping?

past_gossip_rnds = 0; // number of past local gossip rounds

myHopCount = 0; // at what hop did I receive the gossip?

num_rcvdgossips = 0; // number of gossips received so far

initiator = NULL;

Once every protocol period:

if (*gossip_on == TRUE*)

send message (*myHopCount+1, initiator, p*) to gossipTo processes selected uniformly at random from set $G - \text{fromList}$;

past_gossip_rnds++;

if (*past_gossip_rnds > gossipFor*)

gossip_on = FALSE;

On receipt of a gossip message M :

/ received messages processed serially */*

fromList = fromList + { $M.\text{from}$ };

if (*num_rcvdgossips == 0*)

myHopCount = $M.\text{hopNumber}$;

initiator = $M.\text{initiator}$;

gossip_on = TRUE;

num_rcvdgossips++;

if (*num_rcvdgossips > gossipUntil*)

gossip_on = FALSE;

Figure 1. *The Hops Sampling algorithm for estimation. The algorithm is for a one-shot estimation by a single initiator. Continuous estimation requires each gossip message to carry an extra unique identifier for the run number.*

In order to enable estimation, each of the above gossip messages also carries an integer field `hopNumber`, which indicates the number of nodes the message has traversed since the initiator. Before forwarding a gossip, the process stores this hop count in a local variable called `myHopCount`. For multiple received gossip messages, the lowest received `hopNumber` value is remembered. All outgoing gossip messages have their `hopNumber` set to `(myHopCount+1)`.

The initiating process waits for `gossipResult` rounds to elapse before sampling `gossipSample` other processes selected uniformly at random from its membership list. Each sampled process replies with its `myHopCount` value. The average of these values is returned by the algorithm as an estimate of $\log(N)$, where the logarithm's base depends on the parameter settings.

Instead of the initiator sampling the group size, an alternative sampling technique we use in our implementation has the gossip recipients themselves send their hop count values back to the initiator. However, in order to reduce message implosion, the initiating message contains a fixed value `minHopsReporting` specified by the initiator. When a process stops gossiping (if it ever gossips), it sends its `myHopCount` automatically to the initiator (i) with probability 1 if `myHopCount < minHopsReporting`, and (ii) with probability $(1/\text{gossipTo}^{\text{myHopCount}-\text{minHopsReporting}})$ otherwise. Thus, only a small fraction of all hop count values will be received.

Although our protocol description and analysis assume clocks are synchronized to protocol periods, our implementation (Section VI) eliminates this requirement by setting at each process the number of gossip rounds `gossipFor` to the value 1. This means each process gossips only once (to `gossipTo` targets), and this happens as soon as it receives the initiating message. A second modification in the implementation is that the initiator waits for a large, fixed time interval (usually several minutes) before sampling the group for hop counts – this eliminates the need to set `gossipResult`.

III. A. Analysis of the Hops Sampling Algorithm

Theorem 1: Consider a group where the local membership list maintained at each process is $O(\log(N))$, with each entry selected uniformly at random from across the group. Then, the expected time (after gossip initiation) at which a process receives the gossip varies as $\Theta(\log(N))$.

Proof: At the core of the Hops Sampling algorithm is a push gossip protocol, the properties of which are well studied. Bailey [6] and Eugster et al [7] showed the following results for push gossips in a group where processes have partial membership lists chosen uniformly at random from across the group: (i) If the quantity $(\text{gossipTo} * \text{gossipFor})$ is $O(\log(N))$, the gossip spreads to all processes with high probability; (ii) the expected number of rounds for the gossip spread to complete is $\Theta(\log(N))$; (iii) gossip spread in a group with complete membership lists is the same as in a group with partial membership lists of size $O(\log(N))$.

For the purpose of our protocol, we are interested in the *average* time that it takes for the gossip to reach a given process in the group. Clearly (ii) above says it can be no larger than $\Theta(\log(N))$. For the lower bound,

notice that, during each round, the total number of gossip recipients cannot grow by a factor larger than gossipTo . This is a branching process with degree gossipTo and the height of the tree so generated is $O(\log(N))$. This is the lower bound. Thus, the average measured hop count in our algorithm is $\Theta(\log(N))$.

We detail how to determine the value of the logarithm base in the average gossip latency when discussing experiments in Sections V and VI.

III. B. Continuous Version of the Hops Sampling Algorithm

A long-running peer-to-peer application may need to continuously monitor the variation of group size. The continuous variant of the Hops Sampling algorithm can provide this service. It works as follows. The initiator periodically initiates a new one-shot protocol run, each with a unique run identifier. A parameter gossipsAccounted gives the number of recent one-shot estimates to be considered for the continuous estimate. Of these estimates, gossipsDropped of the highest and gossipsDropped of the lowest estimates are dropped before taking the mean of the remaining estimates. The network traffic in this continuous version is bounded since each run terminates w.h.p. in a logarithmic number of rounds.

III. C. Caveats

As pointed out above, the Hops Sampling approach requires a complementary membership protocol that maintains a local membership list at each process. This requirement is easy to satisfy since all distributed systems that are dynamic already run some kind of membership protocol orthogonally. However, only membership protocols that maintain either a complete membership list or a partial membership list with entries chosen uniformly at random from across the group can be used. For example, the gossip-style membership service proposed by van Renesse et al [13] would suffice for such a purpose.

IV. PASSIVE APPROACH - THE INTERVAL DENSITY APPROACH

While the Hops Sampling approach was an *active* probing approach to group size estimation, the Interval Density approach is *passive*. The Interval Density approach is also superior to the Hops Sampling Approach since it requires the complementary membership protocol to only maintain a connected graph, without the need for the membership lists to be chosen uniformly at random.

This approach works by measuring the density of the process identifier space, i.e., the number of processes that have (unique) identifiers lying within an interval of this space. As a first cut, directly sampling the space of process identifiers (e.g., IP addresses + port number) may prove to be inaccurate. For example, if the group were located on a small collection of subnets, the process identifiers would be correlated, and there would be a greater likelihood of error in the group size estimate.

A second approach is to randomize the process identifiers by using a good hash function to map each process' identifier to a point in the real interval $[0,1]$. Cryptographic hash functions such as SHA-1 [8] (or MD-5) can be used: these take as input arbitrary length binary strings, and output a hash of length 160 bits (or 128 bits respectively for MD-5). The hashes can be normalized by dividing with $2^{160} - 1$ (or $2^{128} - 1$ respectively). This is convenient to do in peer-to-peer routing substrates such as Pastry [14] and Chord [15], where virtual "nodeIDs" are already assigned to processes by hashing their identifiers using SHA-1 or MD-5.

The Interval Density approach requires the "initiator" process to passively collect information about the process identifiers that lie in an interval I that is a subset of the interval $[0,1]$. Once this information has been collected, suppose X is the actual number of processes that the initiator finds falling in the interval. Then the estimate for the group size is simply returned as X/I . The passive collection of such process identifier information can be achieved by snooping on the complementary membership protocol running in the overlay – we will describe details of this.

Notice that this lends itself easily to both a one-shot run and a continuous run – the latter can be provided by simply maintaining information about the processes (that lie in the interval I) over a long period of time. Multiple initiators can be supported by selecting the interval specified by the initiator with the lowest identifier, and participating in only that initiator's protocol run.

Below, we first analyze the accuracy of the estimate, then describe adaptive variants of the protocol, and finally discuss the snooping on the membership protocol.

IV. A. Analysis of the Interval Density Approach

Theorem 2: With a good (i.e., uniform) hash function, an interval size I that varies as $O(\log(N)/(\delta^2 \times N))$ suffices to obtain an estimate that is accurate within a factor of δ with high probability, where N be the actual number of non-faulty processes in the group.

Proof: The expected number of processes that fall in the interval I would be $I * N$, where the size of the interval is notated simply as I . Call δ the *accuracy* of the protocol, defined as the multiplicative error in the group size. The likelihood that the estimate for the group size is off by a factor of at most 2δ from the mean is:

$$\Pr[|X/I - N| \leq \delta N] = 1 - \Pr[X/I < (1 - \delta)N] - \Pr[X/I > (1 + \delta)N]$$

If $\delta < 2e - 1$, Chernoff bounds can be used on the latter terms, giving:

$$1 - \Pr[X < I(1 - \delta)N] - \Pr[X > I(1 + \delta)N] \geq 1 - e^{-I N \delta^2/2} - e^{-I N \delta^2/4}$$

If $I > c \log(N)/N$, then the probability that the estimate is accurate would be bounded from below by:

$$(1 - e^{-c \log(N) \delta^2/2} + e^{-c \log(N) \delta^2/4}) \cong 1 - \frac{1}{N^{c \delta^2/2}} - \frac{1}{N^{c \delta^2/4}}$$

Choosing c larger than $1/\delta^2$ suffices to reduce this error asymptotically to zero as N is increased. Thus, to obtain an estimate that is accurate within a constant factor (i.e., δ is a constant) w.h.p., it suffices if I is of length $O(\log(N)/(\delta^2 \times N))$.

Better accuracy can be obtained by using larger interval sizes. If I were of length $O(\sqrt{N}/N)$, the above analysis gives us that the estimate is accurate with probability $(1 - e^{-c \sqrt{N} \delta^2/2} + e^{-c \sqrt{N} \delta^2/4})$. This value asymptotically approaches 1 if $\delta = \sqrt{\log(N)/N^{1/4}}$, which gives a better accuracy than by using a logarithmic interval size.

Practically, of course, it is difficult to set the size of the interval as $O(\sqrt{N}/N)$, without a prior estimate of N . One alternative could be the following. Since the value of (\sqrt{N}/N) decreases as N is increased, if a lower

bound for N is known, the interval can be chosen to be large enough to guarantee a required level of reliability. However, for larger N , a very large number of processes fall inside this interval, thus requiring memory usage at the initiator to grow linearly with group size. Alternatively, this drawback can be addressed by using strategies that adaptively set the interval size. These adaptive strategies also adjust to a bad choice for the location of the interval. They are described next.

The discussion in the remainder of this section applies to only the continuous flavor of the Interval Density approach.

IV. B. *Adapting the Size of the Interval*

An interval I is defined by its *center point* and its *size*. The initiator needs to select a center point for the interval, but the size of this interval itself can be *implicit*, rather than explicit. More specifically, the interval size is determined by remembering a *number* of process identifiers that hash close to the center point. This number is initially set to a small value, but is increased over successive runs until (a) either a predefined threshold of MAXMEMORY processes is reached, or (b) the estimates of group size obtained for successive values of interval length are within a small fraction (e.g., 5%) of each other. Notice that the passive nature of the protocol means that the initiator only selects the interval, but does not communicate it to any other process.

IV. C. *Adapting Interval Location, and Using Multiple Estimation Runs for better Accuracy*

Choosing a good center point for the interval can affect the accuracy of the estimate, especially if process identifiers hash in a rather non-uniform manner into the interval $[0,1]$. Although process identifiers are initially randomized using a cryptographic hash function, this behavior may manifest over a long term due to process arrivals and departures. We describe three approaches for choosing a good interval center:

1. *Random selection.* This is the simplest approach where the initiator randomly selects the center point from the interval $[0,1]$, and then uses this value for all subsequent estimation runs.
2. *Periodically changing selection.* The initiator periodically (e.g., after a few runs) changes the center point of the interval. This could be done either independent of, or dependent on, previously chosen center points.

We call the latter method *self-adjusting interval selection*.

A good self-adjusting interval selection approach is the following. At every stage, history data for the results of a few recent estimation runs is maintained (say the recent 8 runs), along with the intervals that were used for them. The average estimate returned by these recent runs is calculated, and the center point of the run that returned an estimate closest to this average, is used in the next run.

Thus, the center of the interval would continuously move right and left in the $[0,1]$ interval after each round depending on the results of the estimations made on the last rounds. This algorithm will work well if the group size does not change dramatically between two successive estimation runs.

3. *Multiple selections and estimation by mean value.* Instead of changing the center of the interval over time, each estimation run includes multiple runs, each with the interval centered at a different point in $[0,1]$. The mean value of the multiple results is then returned as the estimate.

Our experiments in Section V compare the relative performance of the above three approaches.

IV. D. Snooping on a Membership Maintenance Protocol

We are yet to describe how the initiator passively collects information about process identifiers that lie inside the selected interval I . This is done by snooping on the membership protocol. Below, we describe how this snooping operates on a well-known class of membership protocols called heartbeat-style membership.

In a heartbeat-style membership protocol, e.g., [13], each process p periodically multicasts a heartbeat message (incremented sequence numbers) to *all* other processes in the group. These heartbeats are used to proactively learn about new prospective neighbors in the overlay, as well as for failure detection. The latter is achieved by timing out on the time since the last heartbeat was received by process p from a neighbor process q – this results in p deleting q from its membership list.

At the initiator, the received heartbeat messages can be used to continuously learn about previously unknown processes whose hashed identifiers lie within the interval I being used for estimation, as well as to drop known processes (lying within interval I) that have not responded with an updated heartbeat for a while. Hence, this keeps a running estimate of the number of process identifiers lying within I , and estimates the group size.

Heartbeat-style membership protocols also come in different flavors. We describe in detail how to incorporate the Interval Density approach into the popular gossip-style membership protocol detailed by van Renesse et al [13]. The basic gossip-style membership protocol has each process p periodically (a) increment its own heartbeat counter; (b) select some of its neighbors (defined by the membership list at p), and send to each of these a membership gossip message containing its entire membership list, along with heartbeats. Each process receiving the message merges heartbeat values in the received message with its own membership list. In the modified version of this protocol, used by the Interval Density approach for snooping, a process receiving such a gossip message additionally hashes each previously unknown process identifier appearing in the message, and remembers it only if it lies in the interval I . Notice that the above incorporation is non-intrusive, i.e., it does not affect the normal working of the heartbeat-style membership protocol itself.

The time required for the algorithm to complete is $O(\log(N))$ if gossip messages are allowed to carry up to N heartbeats and identifiers. If the gossip message length is restricted, the time complexity is $O(IN \log N)$.

A final note that applies to both algorithms discussed – a membership protocol that keeps track of the entire list of peers in the system could itself be used to estimate group size. However, most membership protocols for large peer-to-peer overlays do not maintain an entire list due to memory usage and bandwidth constraints. Both presented algorithms require only partial membership lists; the Interval Density algorithm imposes an additional memory requirement to remember some process identifiers.

V. EXPERIMENTAL RESULTS

We evaluate the performance of the Interval Density and the Hops Sampling approaches under two simulation scenarios: (i) Micro-benchmarks, with a static group of a fixed size, and (ii) Trace-based experiments, using trace-log data from Overnet file-sharing network [4] to simulate a dynamic and open group. Our discussion replaces the previous notation “process” with “node”.

The traces from the Overnet network [4] measure the availability of nodes in a 3,000-sized subset of hosts in the deployed Overnet peer-to-peer system. In these traces, the number of hosts that are present in the system

changes by as much as 10% - 25% every hour. We first present individual studies for each of the Hops Sampling algorithm and Interval Density approach, and then compare them against each other.

V. A. Hops Sampling Approach

We assume the membership list at each node to be sampled uniformly at random from the entire set of nodes; lists may contain faulty nodes.

Micro-benchmarks. With parameter settings of `gossipTo=2` and `gossipFor=1`, Figure 2 plots the average number of hops measured in a static network with N nodes versus $\log_2 N$, for different group sizes. The straight line shown in the figure is given by $\log_2 N = 0.9895 * \text{averageNumberOfHops} + 1.3$, and appears to hold for group sizes up to 30,000. The estimated group size is thus calculated as $\text{SizeEst} = 2^{(0.9895 * \text{avHops} + 1.3)}$.

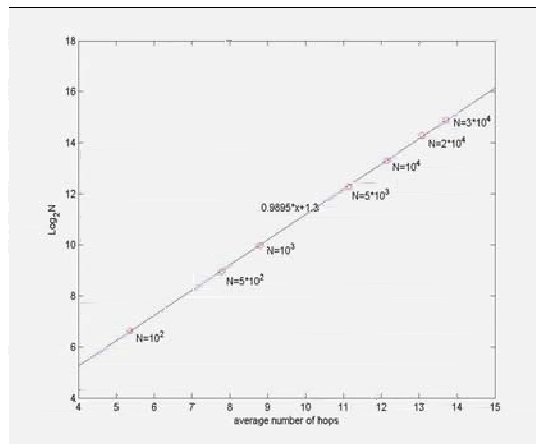


Figure 2. Average number of hops versus $\log_2 N$ for groups with sizes of 100-30,000 nodes.

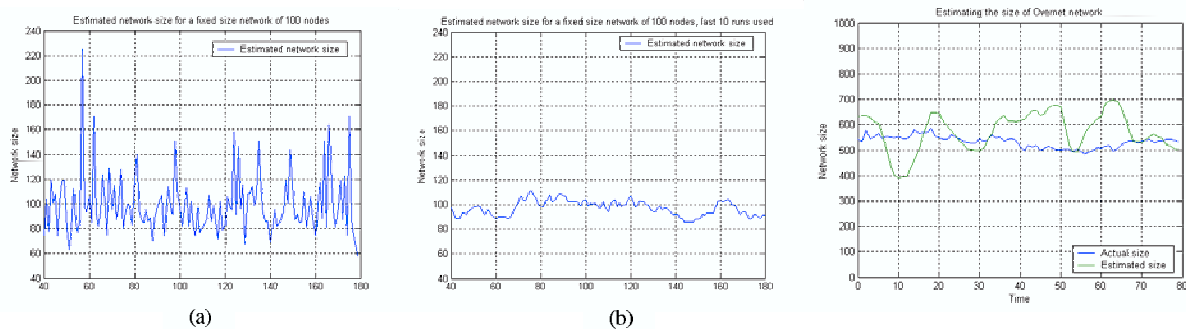


Figure 3. Size Estimation for a static sized network with 100 nodes.

Figure 4. Size Estimation for Overnet network.

Figure 3(a) shows the variation, over time, of the estimated group size, for a static group of size 100. Figure 3(b) shows the effect of averaging over several estimation runs; the continuous protocol version parameters

used are `gossipsAccounted=10` and `gossipsDropped=2`. The benefit of the latter is clear from the figures; the estimate in Figure 3(b) is mostly within 10% of the actual group size.

Trace-based Simulations. The hourly Overnet traces are injected into the simulator at time intervals of 40 timeslots. By “injection”, we mean that the status of the system nodes is updated (as “up” or “down”) at the timeslot we use an Overnet trace. Each time unit of the above plot corresponds to 40 timeslots. The continuous protocol is used, with parameters `gossipTo=2`, `gossipFor=1`, `gossipsAccounted=10` and `gossipsDropped=2`. Figure 4 shows the variation, over time, of the estimated group size. A close look reveals that although the estimate is off by a constant factor (due to a smaller exponent value), it appears to follow and mirror, over time, the variation of the actual system size.

V. B. Interval Density Approach

Interval Density messages are piggybacked on a gossip-style membership protocol, as described earlier in this paper. Similar to [13], entries for processes lying within the interval time out and are marked for deletion if updated heartbeats have not been received for T_{down} time units. Marked entries are deleted after another T_{down} time units; $T_{\text{down}} = 10$ rounds and `MAXMEMORY = 60` for simulations presented in this section.

Micro-benchmark: Static Sized Groups. Figure 5 shows the variation, over time, of the estimate in a static sized group with 10,000 nodes. The random selection approach is applied to decide on the interval used for estimation. The plot shows that the estimate does not converge. This is due to the early expiry of some process entries in the modified gossip-based heartbeat protocol. Figure 6 shows the corresponding timeline when $T_{\text{down}} = \text{infinity}$, so that membership entries never expire. The convergence is within a factor of 5% of the actual group size. Figures 5-6 reveal an interesting property of the Interval Density approach. When membership timeouts (T_{down}) are finite, the estimate is likely to be below the real group size, since some membership entries that lie in the interval are not considered (as they expire). However, at large or infinite values of T_{down} , the estimate can be offset by the local density of hashed node identifiers in the chosen interval. The random approach has 5% accuracy; for better accuracy, the “multiple selections estimation by mean value” approach should be used.

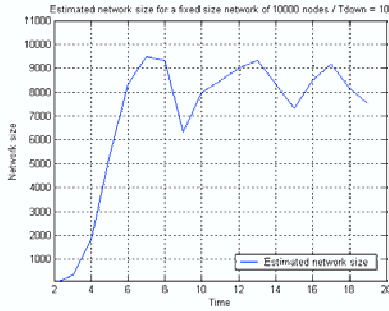


Figure 5. Size estimation for a static sized network of 10,000 nodes, $T_{down} = 10$.

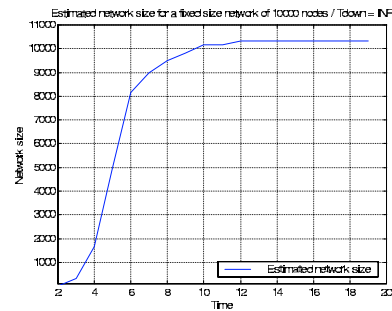


Figure 6. Size estimation for a static sized network of 10,000 nodes, $T_{down} = \text{infinity}$.

Trace-based Simulations: Effect of Adaptive Mechanisms. Figures 7-9 show the behavior, over Overnet traces, of the three adaptive approaches for the selection of interval, presented in section IV.C. The “self-adjusting interval selection” (Figure 8) and the “multiple selections estimation by mean value” (Figure 9) appear to perform the best, although all estimates lie within 20% of the actual group size. As expected, “multiple selections estimation by mean value” has the least variance among all adaptive methods since multiple intervals are sampled, and an average is taken over them.

Highly Dynamic Groups. Although the fraction of nodes joining and leaving the system is as high as 25% in the Overnet traces, the reader would have observed from the previous experiments that the actual group size itself does not vary much. To model a dynamic group, we modify the Overnet traces so that an additional set of arrivals and departures is added to the original ones, causing the number of nodes in the network to vary between time intervals. The self-adjusting method is used for providing estimations. The results are presented in Figure 10, and show good behavior in spite of highly dynamic groups. A look at Figure 10 also shows that the estimate closely follows and mirrors the variation of the actual group size over time.

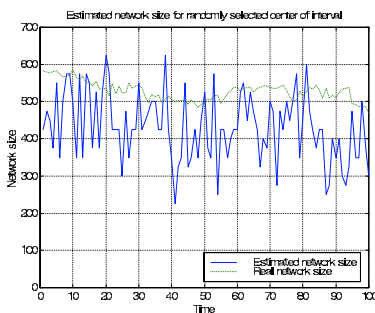


Figure 7. Estimate calculated by the random interval selection approach.

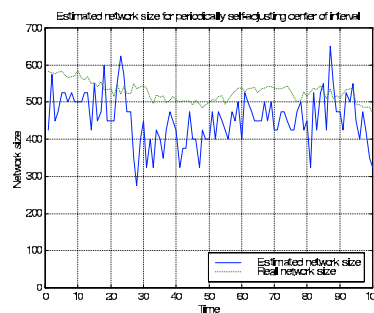


Figure 8. Estimate calculated by the self-adjusting interval selection approach.

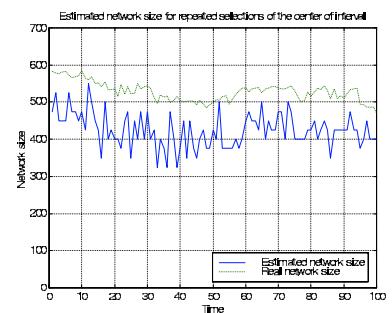


Figure 9. Estimate calculated by the estimation by mean value approach.

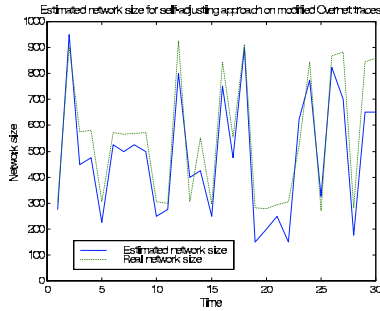


Figure 10. Behavior of self-adjusting interval selection in a highly dynamic group (Overnet traces, with additional random joins and failures).

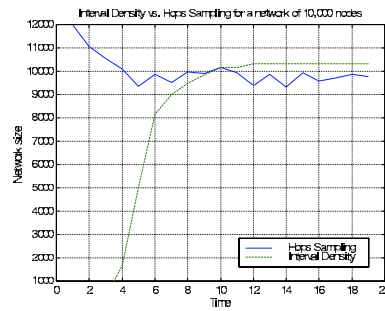


Figure 11. Comparing the accuracy of the Interval Density and the Hops Sampling approaches for a large fixed size network of 10,000 nodes.

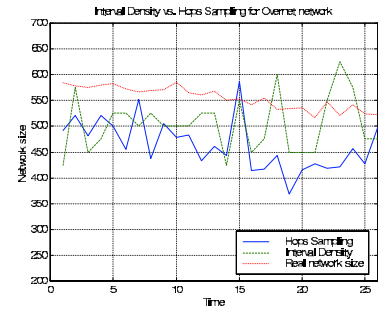


Figure 12. Comparing the accuracy of the Interval Density and the Hops Sampling approaches for Overnet traces.

V. C. Hops Sampling Algorithm versus Interval Density Algorithm

Figures 11 and 12 compare the proposed approaches in static (fixed size) and dynamic (Overnet-trace based) groups. The gossip parameter settings are $\text{gossipTo}=2$, $\text{gossipFor}=1$, $\text{gossipsAccounted}=10$ and $\text{gossipsDropped}=2$. Moreover, $\text{MAXMEMORY} = 60$ and $T_{\text{down}} = \text{infinity}$ in Figure 11. The self-adjusting interval selection method, with $\text{MAXMEMORY} = 120$ and $T_{\text{down}} = 10$, is used for Figure 12.

For static sized groups, both the Interval Density and the Hops Sampling approaches converge quickly to within 5% of the actual group size (Figure 11), with the Hops Sampling approach performing slightly better. For dynamic behavior due to Overnet traces, the estimate from the Interval Density approach is slightly better than that from the Hops Sampling approach (Figure 12). However, the Interval Density approach has the advantage that it does not require membership entries to be chosen uniformly at random from across the group.

VI. PEERCOUNTER: A SIZE ESTIMATION UTILITY

In this section, we present PeerCounter. PeerCounter is a generic module containing both algorithms present in this paper, and it can be incorporated into an arbitrary peer-to-peer overlay. We first briefly outline PeerCounter’s Application Programming Interface (API) and then demonstrate evaluation results obtained by applying this prototype implementation of our estimation schemes on a cluster of 33 Linux servers located in the University of Illinois. For our evaluation analysis we use a topology generator script in order to set up membership lists used for gossiping. Any synchronization assumptions made in our analysis or simulations are relaxed here since PeerCounter does not require any synchronization among processes.

VI. A. PeerCounter's API

PeerCounter is a command line application implemented in Java. It takes as parameters the server port where to run the application, a symbol that indicates which of the two algorithms is being used for estimation, the Hops Sampling or the Interval Density, and the membership list of the calling process. The membership list can be generated using a topology generator utility or, in case PeerCounter is incorporated into a standard overlay, this may be provided by the underlying application. PeerCounter's API can be used to let applications layered under it use its estimation schemes. PeerCounter exports the following operations:

alg = new HopCountingAlgorithm(members) or alg = new IntervalCountingAlgorithm(members)	Instantiates one of the estimation algorithms <i>members</i> holds the membership list of the calling process
alg.startCounting()	Causes the gossip stream to start spreading over the group
alg.getSize()	Causes a size estimate to be returned

Source code of the PeerCounter implementation can be found at: http://netfiles.uiuc.edu/psaltoul/www/peerCounter/index_files/download.htm

VI. B. Performance Results

The Hops Sampling Scheme. The parameters used here are `gossipTo=2`, `gossipFor=1`, `minHopsReporting=4`, `gossipsAccounted=5` and `gossipsDropped=1`. We test both the standard estimation scheme, which uses the average hop distance from processes sampled by the initiator (as presented in Figure 1), as well as the alternative sampling scheme, mentioned in section III, where hop count values reported back at the initiator by processes with certain `myHopCount` values, are taken into account for size estimation. The equation used for the first case is $size = 2^{(averageNumberOfHops+1.75)}$. This is very close to the one extracted by using our simulations in Section V, with only a small difference in the constant factor. Figure 13 demonstrates the estimates provided by applying our PeerCounter application on static groups with sizes of 80, 806 and 6430 nodes. Estimates are found to lie on average within 10% of the actual group size.

The Interval Density Scheme. Figure 14 shows the estimates provided by PeerCounter for static groups of different sizes when the one-shot flavor of the Interval Density approach is being used and a random selection of the interval is being made. The parameters used here are `MAXMEMORY = 100` and `Tdown = infinity`. As shown in Section V, the accuracy of estimations is expected to be better when adaptive methods or multiple selections

and estimation by mean value are applied in order to choose the interval location. Moreover, a continuous running of the application should cause the estimates to converge to some value near the actual group size, when membership timeouts are ignored, as studied earlier in this paper. Figure 15, which shows results for the continuous flavor of the Interval Density approach, indicates that this actually holds for our prototype implementation as well.

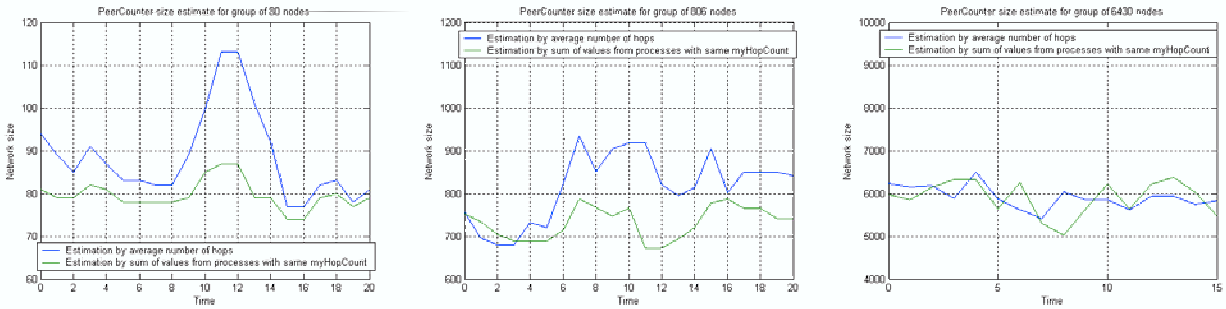


Figure 13. Estimates provided from PeerCounter for using Hops Sampling scheme on groups of 80, 806 and 6430 nodes.

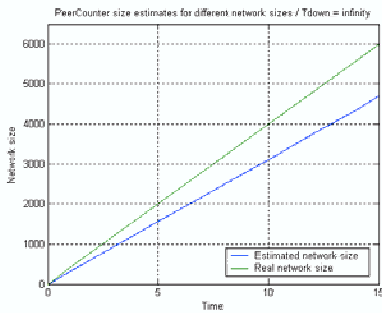


Figure 14. Estimates provided by PeerCounter for one-shot Interval Density scheme with random interval selection.

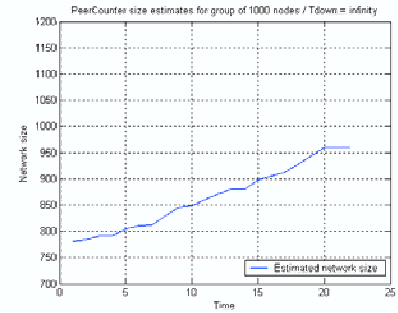
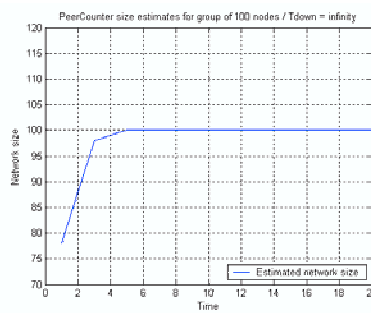


Figure 15. Estimates provided by PeerCounter for groups of 100 and 1,000 nodes when the continuous Interval Density scheme with random interval selection is being used.

VII. CONCLUSIONS

Estimating the size of a decentralized group of processes connected within an overlay (e.g., in a peer-to-peer system, or in the Grid) is a difficult problem, especially when the group is dynamic and contains thousands of processes. Previous solutions to the problem make assumptions that may be unscalable or limit applicability. In this paper, we have proposed two new approaches for estimation – an active approach called Hops Sampling and a passive approach called Interval Density. The only requirements for these algorithms are the existence of complementary membership protocols, which are already present as a part of most overlays.

Micro-benchmarks (for static groups) and trace-based simulations have shown that both above approaches are able to obtain an estimate within a few percentage points of the actual group size. The Hops Sampling approach gives an accurate estimate if the local membership lists at each process are chosen uniformly at random. The Interval Density approach is more practical, since its accuracy only requires the overlay graph among all processes in the system to be connected. We have made these algorithms available as an open-source size estimation utility called PeerCounter that can be incorporated into any peer-to-peer overlay.

REFERENCES

- [1] B. Awerbuch and C. Scheideler, "*Robust Distributed Name Service*", In Proc. 3rd International Workshop on Peer-to-Peer Systems (IPTPS), 2004.
- [2] N. T. J. Bailey, "*Epidemic Theory of Infectious Diseases and its Applications*", Hafner Press, 2nd Edition, 1975.
- [3] M. Bawa, H. Garcia-Molina, A. Gionis and R. Motwani, "*Estimating aggregates on a peer-to-peer network*", Technical Report, Dept. of Computer Science, Stanford University, 2003.
- [4] R. Bhagwan, Overnet availability traces: <http://ramp.ucsd.edu/projects/recall/download.html>.
- [5] K. P. Birman, M. Hayden, O. Ozkasap, Z. Xiao, M. Budiu and Y. Minsky, "*Bimodal Multicast*", ACM Transactions on Computer Systems, Vol. 17, No. 2, pp. 41-88, May 1999.
- [6] A. Demers, D. Greene, C. Hauser, W. Irish, J. Larson, S. Shenker, H. Sturgis, D. Swinehart and D. Terry, "*Epidemic algorithms for replicated database maintenance*", In Proc. Symposium on Principles of Distributed Computing, pp. 1-12, August 1987.
- [7] P. Th. Eugster, R. Guerraoui, S. B. Handurukande, A.M. Kermarrec and P. Kouznetsov, "*Lightweight probabilistic broadcast*", ACM Transactions on Computer Systems, Vol. 21, No. 4, pp. 341-374, November 2003.
- [8] FIPS 180-1, "*Secure Hash Standard*", NIST, US Department of Commerce, Washington D.C., April 1995.
- [9] I. Gupta, K. Birman, P. Linga, A. Demers and R. van Renesse, "*Kelips: building an efficient and stable P2P DHT through increased memory and background overhead*", In Proc. 2nd International Workshop on Peer-to-peer Systems, Springer LNCS 2735, pp. 160-169, 2003.
- [10] R. M. Karp, C. Schindelhauer, S. Shenker and B. Vocking, "*Randomized rumor spreading*", IEEE Symposium on Foundations of Computer Science, pp. 565-574, 2000.
- [11] A.-M. Kermarrec and L. Massoulie, Private communication, 2003.
- [12] D. Malkhi and K. Horowitz, "*Estimating network size from local information*", ACM Information Processing Letters, Vol. 88, Issue 5, pp. 237-243, December 2003.
- [13] R. van Renesse, Y. Minsky and M. Hayden, "*A gossip-style failure detection service*", In Proc. Middleware '98, Lancaster, England, September 1998.
- [14] A. Rowstron and P. Druschel, "*Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer Systems*", In Proc. of IFIP / ACM Middleware, pp. 329-350, November 2001.
- [15] I. Stoica, R. Morris, D. Karger, F. Kaashoek and H. Balakrishnan, "*Chord: A scalable peer-to-peer lookup service for Internet applications*", In Proc. ACM SIGCOMM Conference, August 2001.