



Computer Science and Artificial Intelligence Laboratory  
Technical Report

MIT-CSAIL-TR-2005-052  
MIT-LCS-TR-999

August 11, 2005

---

**Self-Stabilizing Mobile Node Location  
Management and Message**

Shlomi Dolev, Limor Lahiani, Nancy Lynch, Tina Nolte

# Self-Stabilizing Mobile Node Location Management and Message Routing

Shlomi Dolev\*      Limor Lahiani\*      Nancy Lynch†      Tina Nolte†

August 11, 2005

## Abstract

We present simple algorithms for achieving self-stabilizing location management and routing in mobile ad-hoc networks. While mobile clients may be susceptible to corruption and stopping failures, mobile networks are often deployed with a reliable *GPS oracle*, supplying frequent updates of accurate real time and location information to mobile nodes. Information from a GPS oracle provides an external, shared source of consistency for mobile nodes, allowing them to label and timestamp messages, and hence aiding in identification of, and eventual recovery from, corruption and failures. Our algorithms use a GPS oracle.

Our algorithms also take advantage of the *Virtual Stationary Automata* programming abstraction, consisting of mobile clients, virtual timed machines called virtual stationary automata (VSAs), and a local broadcast service connecting VSAs and mobile clients. VSAs are distributed at known locations over the plane, and emulated in a self-stabilizing manner by the mobile nodes in the system. They serve as fault-tolerant building blocks that can interact with mobile clients and each other, and can simplify implementations of services in mobile networks.

We implement three self-stabilizing, fault-tolerant services, each built on the prior services: (1) VSA-to-VSA geographic routing, (2) mobile client location management, and (3) mobile client end-to-end routing. We use a greedy version of the classical depth-first search algorithm to route messages between VSAs in different regions. The mobile client location management service is based on *home locations*: Each client identifier hashes to a set of home locations, regions whose VSAs are periodically updated with the client's location. VSAs maintain this information and answer queries for client locations. Finally, the VSA-to-VSA routing and location management services are used to implement mobile client end-to-end routing.

**Keywords:** virtual infrastructure, location management, home locations, end-to-end routing, hash functions, self-stabilization, GPS oracle

---

\*Department of Computer Science, Ben-Gurion University of the Negev, Beer-Sheva, 84105, Israel. Partially supported by IBM faculty award, NSF grant and the Israeli ministry of defense. Email: {dolev, lahiani}@cs.bgu.ac.il.

†MIT Computer Science and Artificial Intelligence Laboratory, The Stata Center 32-G670, Cambridge, MA 02139, USA. Supported by DARPA contract F33615-01-C-1896, NSF ITR contract CCR-0121277, and USAF, AFRL contract FA9550-04-1-0121. Email: {lynch, tnolte}@theory.csail.mit.edu.

# 1 Introduction

A system with no fixed infrastructure in which mobile clients may wander in the plane and assist each other in forwarding messages is called an ad-hoc network. The task of designing algorithms for constantly changing networks is difficult. Highly dynamic networks, however, are becoming increasingly prevalent, especially in the context of pervasive and ubiquitous computing, and it is therefore important to develop and use techniques that simplify this task. In addition, mobile nodes in these networks may suffer from crash failures or corruption faults, which cause arbitrary changes to their program states. Self-stabilization [4, 5] is the ability to recover from an arbitrarily corrupt state. This property is important in long-lived, chaotic systems where certain events can result in unpredictable faults. For example, transient interference may disrupt wireless communication, violating our assumptions about the broadcast medium.

Mobile networks are often deployed in conjunction with “reliable” GPS services, supplying frequent updates of real time and region information to mobile nodes. While the mobile clients may be susceptible to corruption and stopping failures, the GPS service may not be. Each of our algorithms utilizes such a reliable *GPS oracle*. Information from this oracle provides an external, shared source of consistency for mobile nodes, allowing them to label and timestamp their messages, and hence, aiding in identification of, and recovery from, corruption and stopping failures.

In this paper we describe self-stabilizing algorithms that use a reliable GPS oracle to provide geographic routing, a mobile client location management service, and a mobile client end-to-end routing service. Each service is built on the prior services such that the composition of the services remains self-stabilizing [11]. In order to route location information between geographic regions, we use a greedy version of the classical depth-first search algorithm. This service is then used to help implement the location management service; each mobile client identifier hashes to a set of home locations, geographical regions that are periodically updated with the location of the client, and that are responsible for then answering queries about the client’s location. Both of these services are then used to implement point-to-point routing between mobile clients in the network.

In order to simplify the implementations of the location management and routing services, we mask the unpredictable behavior of mobile nodes by using a self-stabilizing *virtual* infrastructure, consisting of mobile client automata, timing-aware and location-aware machines at fixed locations, called *Virtual Stationary Automata* (VSAs) [8, 9], that mobile clients can interact with and use to coordinate their actions, and a local broadcast service connecting VSAs and mobile clients.

**Self-stabilization and *GPS oracles*.** Traditionally, studies of self-stabilizing systems are concerned with those systems that can be started from arbitrary configurations and eventually regain consistency *without external help*. However, mobile clients often have access to some reliable external information from a service such as GPS. Each of our algorithms in this paper uses an external GPS service (or an equivalent service) as a reliable *GPS oracle*, providing periodic real time clock and location updates, to base stabilization upon; our algorithms use timestamps and location information to tag events. In an arbitrary state, recorded events may have corrupted timestamps. Corrupted timestamps indicating future times can be identified and reset to predefined values; new events receive newer timestamps than any in the arbitrary initial state. This eventually allows nodes in the system to totally order events. We use the eventual total order to provide consistency of information and distinguish between incarnations of activity (such as retransmissions of messages).

**Virtual Stationary Automata programming layer.** In prior work [8, 7, 6], we developed a notion of “virtual nodes” for mobile ad hoc networks. A virtual node is an abstract, relatively well-behaved active node that is implemented using less well-behaved real physical nodes. The GeoQuorums algorithm [7] proposes storing data at fixed locations; however it only supports atomic objects, rather than general automata. A more general virtual mobile automaton is suggested in [6]. Finally, the virtual automata presented in [8, 9] (and used here) are more powerful than those of [6], providing timing capabilities needed for many applications. These automata are stationary and arranged in a connected pattern similar to that of a traditional wired network.

The static infrastructure we use in this paper includes fixed, timed virtual machines with an explicit notion of real time, called *Virtual Stationary Automata* (VSAs), distributed at known locations over the plane [8, 9]. Each VSA represents a predetermined geographic area and has broadcast capabilities similar to those of the

mobile physical nodes, allowing nearby VSAs and mobile nodes to communicate with one another. Many algorithms depend significantly on timing, and it is reasonable to assume that many mobile nodes have access to reasonably synchronized clocks. In the VSA layer, VSAs also have access to *virtual* clocks, guaranteed to not drift too far from real time. The layer provides mobile nodes with a fixed virtual infrastructure, reminiscent of more traditional and better understood wired networks, with which to coordinate their actions.

Our clock-enabled VSA layer is emulated by physical mobile nodes in the network. Each physical node is periodically informed its region by the GPS. A VSA for a particular region is then emulated by a subset of the mobile nodes in its region: the VSA state is maintained in the memory of the physical nodes emulating it, and the physical nodes perform VSA actions on behalf of the VSA. If no physical nodes are in the region, the VSA fails; if physical nodes later arrive, the VSA restarts.

An important property of the VSA layer implementation described in [8, 9] is that it is self-stabilizing. Corruption failures at physical nodes can result in inconsistency in the emulation of a VSA. Our implementation, however, can recover after corruptions to correctly emulate a VSA. To algorithms run on the VSA layer, the VSA simply appears to suffer from a corruption.

**Geographic/ VSA-to-VSA routing.** A basic service running on the VSA layer that we describe and use repeatedly is that of VSA-to-VSA (region-to-region) routing (VtoVComm), providing a form of geocast. GeoCast algorithms [24, 3], GOAFR [19], and algorithms for “routing on a curve” [23] route messages based on the location of the source and destination, using geography to delivery messages efficiently. GPSR [17], AFR [20], GOAFR+ [19], polygonal broadcast [10], and the asymptotically optimal algorithm [20] are algorithms based on greedy geographic routing algorithms, forwarding messages to the neighbor that is geographically closest to the destination. The algorithms also address “local minimum situations”, where the greedy decision cannot be made. GPSR, GOAFR+, and AFR achieve, under reasonable network behavior, a linear order expected cost in the distance between the sender and the receiver. We implement VSA-to-VSA routing using a persistent greedy depth-first search (DFS) routing algorithm that runs on top of the VSA layer’s fixed infrastructure. Our scheme is an application of the classical DFS algorithm in a new setting.

**Location management.** Finding the location of a moving client in an ad-hoc network is difficult, much more so than in cellular mobile networks where a fixed infrastructure of wired support stations exist (as in [16]), or in sensor networks where some approximation of a fixed infrastructure may exist [2]. A *location service* in ad-hoc networks is a service that allows any client to discover the location of any other client using only its identifier. The basic paradigm for location services that we use here is that of a home location service: Hosts called *home location servers* are responsible for storing and maintaining the location of other hosts in the network [1, 14, 21]. Several ways to determine the sets of home location servers, both in the cellular and entirely ad-hoc settings, have been suggested.

The locality aware location service (LLS) in [1] for ad-hoc networks is based on a hierarchy of lattice points for destination nodes, published with locations of associated nodes. Lattice points can be queried for the desired location, with a query traversing a spiral path of lattice nodes increasingly distant from the source until it reaches the destination. Another way of choosing location servers is based on quorums. A set of hosts is chosen to be a *write* quorum for a mobile client and is updated with the client’s location. Another set is chosen to be a *read* quorum and queried for the desired client location. Each *write* and *read* quorum has a nonempty intersection, guaranteeing that if a *read* quorum is queried, the results will include the latest location of the client written to a *write* quorum. In [14], a uniform quorum system is suggested, based on a virtual backbone of quorum representatives. Geographic quorums based on the focal points abstraction are suggested in [7].

Location servers can also be chosen using a hash table. Some papers [21, 15, 25] use geographic locations as a repository for data. These use a hash to associate each piece of data with a region of the network and store the data at certain nodes in the region. This data can then be used for routing or other applications. The Grid location service (GLS) [21] maps client ids to geographic coordinates. A client  $C_p$ ’s location is saved by clients closest to the coordinates  $p$  hashes to.

The location management scheme we present here is based on the hash table concept and built on top of the VSA layer and VSA-to-VSA routing service. VSAs and mobile clients are programmed to form a self-stabilizing, fault-tolerant distributed data structure for location management, where VSAs serve as home locations for mobile clients. Each client’s id hashes to a VSA region, the client’s home location, whose VSA is responsible for maintaining the location of the client. Whenever a client node  $C_p$  would like to locate

<p><b>System constants:</b></p> <p><math>R</math>, a fixed closed connected region of the 2-D plane.  <math>U</math>, a finite set of ids for subregions of <math>R</math>.  <math>m</math>, the size of <math>U</math>.  <math>region</math>, a mapping from <math>U</math> to connected subsets of <math>R</math>.  <math>nbrs</math>, a symmetric relation between ids in <math>U</math>.  <math>r_{virt}</math>, the supremum distance between points in <math>u</math> and <math>v</math> for any regions <math>u, v</math> where <math>u \in nbrs(v)</math>.  <math>P</math>, a finite set of client node ids where <math>P \cap U = \emptyset</math>.  <math>v_{max}</math>, the maximum client node speed.</p>	<p><math>\epsilon_{sample}</math>, the GPS sample period.  <math>d</math>, the broadcast message delay.  <math>e</math>, the delay factor for VSA outputs.  <math>tll_{VtoV} &gt; d</math>, the VtoVComm message delay.  <math>tVSAcor</math>, the VSA stabilization time.</p> <p><b>System variables:</b></p> <p><math>now \in \mathbb{R}</math>, a clock variable, representing real time.  <math>loc</math>, a continuously updated array of locations in <math>\mathbb{R}</math> of mobile nodes, indexed by node id.</p>
<p>Figure 1: System constants and variables.</p>	

another client node  $C_q$ ,  $C_p$  would compute the home location of  $C_q$  by applying a predefined global hash function to  $C_q$ 's id, and query the region represented by the result of that hash for  $C_q$ 's location. In order for our scheme to tolerate crash failures of a limited number of VSAs, each mobile client id actually maps to a set of VSA home locations; the hash function returns a sequence of region ids as the home locations. We can use any hash function that provides a sequence of region identifiers; one possibility is a *permutation hash function*, where permutations of region ids are lexicographically ordered and indexed by client id.

**End-to-end routing.** Another basic, but difficult to provide, service in mobile networks is end-to-end routing. Our self-stabilizing implementation of a mobile client end-to-end communication service is simple, given VSA-to-VSA routing and the home location service. A client sends a message to another client by using the home location service to discover the destination client's region and then has a local VSA forward the message to the region using the VSA-to-VSA service.

**Paper organization.** The rest of the paper is organized as follows: The system model and the virtual automata layer are described in the next section. In Section 3 we describe the problem specifications we are interested in. Section 4 describes the VSA-to-VSA communication implementation. In Section 5 we describe the implementation of the home location service. In Section 6 we present the implementation of the end-to-end routing service. Concluding remarks appear in Section 7.

## 2 Datatypes and system model

The system consists of a 2-D bounded region plane, where broadcast-enabled, GPS-updated mobile client nodes are deployed. We assume the *Virtual Stationary Automata* programming abstraction [8], which includes both the mobile client nodes and virtual stationary automata (VSAs) the real nodes emulate, as well as a local broadcast service, V-bcast, between them (see Figure 2). In this section we formally describe the system, including: (1) the network tiling, (2) the model for the GPS-augmented mobile clients deployed in the network, (3) the model for the virtual nodes deployed in the network, and (4) the specification for the local broadcast service in the network. A summary table of datatypes, constants, and variables is in Figure 1.

### 2.1 Network tiling

The deployment space of the network is assumed to be a fixed, closed, and bounded connected region of the 2-D plane called  $R$ .  $R$  is partitioned into known connected subregions called *regions*, with unique ids drawn from the set of region identifiers  $U$ . In practice it may be convenient to restrict regions to be regular polygons such as squares or hexagons. We define a neighbor relation  $nbrs$  on ids from  $U$ . This relation holds for any two region identifiers  $u$  and  $v$  where the supremum distance between points in  $u$  and  $v$  is bounded by a constant  $r_{virt}$ .

### 2.2 Client nodes

For each  $p$  in the set of physical node identifiers  $P$ , we assume a mobile timed I/O automaton client  $C_p$ , whose location in  $R$  at any time is referred to as  $loc(p)$ . Mobile client speed is bounded by a constant  $v_{max}$ . Clients receive region and time information from the GPS oracle. A  $GPSupdate(u, now)_p$  happens every  $\epsilon_{sample}$  time at each client  $C_p$ , indicating to the client the region  $u$  where it is currently located and the

current time  $now$ . Clients accept this  $now$  real-time clock variable as the value of their own local clock. For simplicity, this local variable progresses at the rate of real time. This implies that, outside of failures, the local value of  $now$  will equal real time.

Each client  $C_p$  is equipped with a local broadcast service V-bcast (see Section 2.4), allowing it to communicate with its and neighboring regions' VSAs and clients with  $\text{bcast}(m)_p$  and  $\text{brcv}(m)_p$ .

Clients are susceptible to stopping and corruption failures. After a stopping failure, a client performs no additional local steps until restarted. If restarted, it starts again from an initial state. If a node suffers from a corruption, it experiences a nondeterministic change to its program state.

Additional arbitrary external interface actions and local state used by algorithms running at the client are allowed. For simplicity local steps are assumed to take no time.

### 2.3 Virtual Stationary Automata (VSAs)

Here we describe VSAs; a self-stabilizing implementation of such machines using a GPS oracle and the physical mobile nodes in the system can be found in [8, 9]. An abstract VSA is a timing-enabled virtual machine that may be emulated by the physical mobile nodes in its region in the network. We formally describe a timed machine for region  $u$ ,  $V_u$ , as a TIOA whose program is a tuple of its action signature,  $\text{sig}_u$ , valid states,  $\text{states}_u$ , a start state function mapping clock values to start states,  $\text{start}_u$ , a discrete transition function,  $\delta_u$ , and a set of valid trajectories,  $\tau_u$ . Trajectories [18] describe state evolution over intervals of time. The state of  $V_u$  is referred to collectively as  $vstate$  and is assumed to include a variable corresponding to real time,  $vstate.now$ .

To guarantee that we can emulate a VSA using physical mobile nodes, its interface must be emulatable by the nodes. Hence, a VSA  $V_u$ 's external interface is restricted to be similar, including only stopping failure, corruption, and restart inputs, and the ability to broadcast and receive messages. Corruption failures result in a nondeterministic change to  $vstate$ .

Since a VSA is emulated by physical nodes (corresponding to clients) in its region, its failures are defined in terms of client failures in its region: (1) If no clients are in the region, the VSA is crashed, (2) If no failures of clients (corruption or stopping) occurs in an alive VSA's region over some interval, the VSA does not suffer a failure during that interval, and (3) A VSA may suffer a corruption only if a mobile client in its region suffers a corruption; the self-stabilizing implementation of a VSA in [8, 9] guarantees that within  $t_{VSAcor}$  of an arbitrary configuration of the emulation, the emulation's external trace will look like that of the abstract VSA, starting from a corrupted abstract state.

While an emulation of  $V_u$  would ideally be identical to a legitimate execution of  $V_u$ , an abstraction must reflect that, due to message delays or node failure, the emulation might be behind real time, appearing to be delayed in performing outputs by up to some time  $e$ . The emulation is then a *delay-augmented TIOA*, an augmentation of  $V_u$  with timing perturbations, represented with buffers  $\text{Dout}[e]_u$ , composed with  $V_u$ 's outputs. The buffer delays messages by a nondeterministic time  $[0, e]$ , where  $e$  is more than V-bcast's broadcast delay,  $d$  (see Section 2.4). Programs must take into account  $e$ , as they do  $d$ .

### 2.4 Local broadcast service (V-bcast)

Communication is in the form of local broadcast service V-bcast, with broadcast radius  $r_{virt}$  and message delay  $d$ . It allows communication between VSAs and clients in the same or neighboring regions. The service allows the broadcasting and receiving of message  $m$  at each port  $i \in P \cup U$  through  $\text{bcast}(m)_i$  and  $\text{brcv}(m)_i$ .

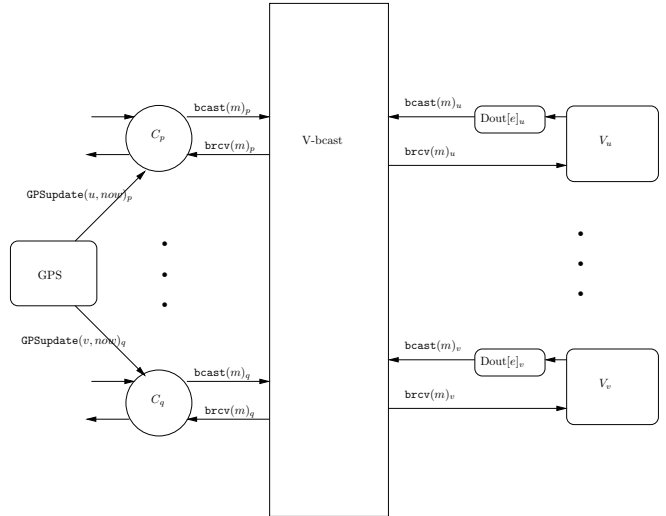


Figure 2: Virtual Stationary Automata layer. VSAs and clients communicate locally using V-bcast. VSA outputs may be delayed in  $\text{Dout}$ .

We assume that V-bcast guarantees two properties between VSAs and between VSAs and clients: integrity and reliable local delivery. *Integrity* guarantees that for any  $\text{brvc}(m)_i$  that occurs, a  $\text{bcast}(m)_j$ ,  $j \in P \cup U$  previously occurred. *Reliable local delivery* roughly guarantees that a transmission will be received by nearby ports: If port  $i$ , where  $i$  is a client or VSA port in any region  $u$ , transmits a message, then every port  $j$ , whether a client or VSA port, in region  $u$  or neighboring regions during the entire time interval starting at transmission and ending  $d$  later receives the message by the end of the interval. (For this definition, due to GPSupdate lag, a client is still said to be “in” region  $u$  even if it has just left region  $u$  but has not yet received a GPSupdate with the change.)

In practice, a broadcast service has bounded buffers. We assume buffers are large enough that overflows do not occur in normal operation. In the event of overflow, overflow messages are lost.

### 3 Problem specifications

We describe the services we will build over the VSA layer: VSA-to-VSA routing, a location service, and client-to-client routing, and describe our requirement that implementations be self-stabilizing.

The following constants (explained/used shortly) are globally known: (1)  $f < m$ , a limit on “home location” VSA failures for a client, (2)  $h$ , a function mapping each client id to a sequence of  $f + 1$  distinct region ids, (3)  $\text{ttl}_{VtoV} > d$ , delivery time for the VtoVComm service, (4)  $\text{ttl}_{HLS} \geq \epsilon_{\text{sample}} + 2d + 3e + 2\text{ttl}_{VtoV}$ , response time of the location management service, and (5)  $\text{ttl}_{hb}$ , a refresh period. We assume the following client mobility and VSA crash failure conditions:

- (1) Each client spends at least  $\epsilon_{\text{sample}}$  time in a region before moving to another region,
- (2) At any time, each alive client’s current region or a neighboring region has a non-crashed VSA that remains alive for an additional  $\text{ttl}_{HLS}$  time,
- (3) For any interval of length  $\text{ttl}_{VtoV} + e$ , two VSAs alive over the interval are connected via at least one path of non-crashed VSAs over the entire interval, and
- (4) For any interval of length  $\text{ttl}_{hb} + 2\text{ttl}_{VtoV} + 2e + d$ , and any alive client  $q$ , at least one VSA from  $h(q)$  does not crash during the interval.

#### 3.1 VSA-to-VSA communication service (VtoVComm) specification

The first service is an inter-VSA routing service, where a VSA from some region  $u$  can send a message  $m$  through  $\text{VtoVsend}(v, m)_u$  to a VSA in another (potentially non-neighboring) region  $v$ . Region  $v$ ’s VSA later receives  $m$  through  $\text{VtoVrcv}(m)_v$ . The service guarantees two properties:

- (1) If a VSA at region  $u$  performs a  $\text{VtoVsend}(v, m)$ , and both region  $u$  and  $v$  VSAs are alive over the time interval beginning with the send and ending  $\text{ttl}_{VtoV}$  time later, then the VSA at region  $v$  performs a  $\text{VtoVrcv}(m)$  before the end of the interval, and
- (2) If a message is received at some VSA, it was previously sent to that VSA.

#### 3.2 Location service specification

A location service answers queries from clients for the locations of other clients. A client node  $p$  can submit a query for a recent region of client node  $q$  via a  $\text{HLquery}(q)_p$  action. If few home location failures occur and  $q$  has been in the system for a sufficient amount of time, the service responds within bounded time with a recent region location of  $q$ ,  $qreg$ , through a  $\text{HLreply}(q, qreg)_i$  action.

To be more exact, the location service guarantees that if a client  $p$  performs a  $\text{HLquery}$  to find an alive client  $q$  that has been in the system longer than  $\epsilon_{\text{sample}} + d + \text{ttl}_{VtoV} + e + \text{ttl}_{HLS}$  time, and client  $p$  does not crash or move to a different region for  $\text{ttl}_{HLS}$  time, then:

- (1) Within  $\text{ttl}_{HLS}$  time, client  $p$  will perform a  $\text{HLreply}$  with a region for  $q$ , and
- (2) If  $p$  performs a  $\text{HLreply}(q, qreg)$ , then  $p$  had requested  $q$ ’s location and  $q$  was either: (a) alive in region  $qreg$  within the last  $\text{ttl}_{HLS}$  time, or (b) failed for at most  $\text{ttl}_{hb} + \text{ttl}_{HLS} - \epsilon_{\text{sample}}$  time.

#### 3.3 Client end-to-end routing (EtoEComm) specification

End-to-end routing is an important application for ad-hoc networks. The V-bcast service provides a local broadcast service where VSAs and clients can communicate with VSAs and clients in neighboring regions. VtoVComm allows arbitrary VSAs to communicate. End-to-end routing (EtoEComm) allows arbitrary

clients to communicate: a client  $p$  sends message  $m$  to client  $q$  using  $\text{send}(q, m)_p$ , which is received by  $q$  in bounded time via  $\text{receive}(m)_q$ .

If clients  $p$  and  $q$  do not crash for  $t_{HLS}$  time, clients do not change regions for  $t_{HLS}$  time after a send, and  $q$  has been in the system at least  $t_{HLS} + \epsilon_{\text{sample}} + d + t_{VtoV} + e$  time, then:

- (1) If client  $p$  sends message  $m$  to  $q$ ,  $q$  will receive  $m$  within  $t_{HLS} + 2d + 2e + t_{VtoV}$  time, and
- (2) Any message received by a client was previously sent to the client.

### 3.4 Self-stabilizing implementations

We require implementations of the above services to be self-stabilizing. A system configuration is *safe* with respect to a specification and implementation if any admissible execution fragment of the implementation starting from the configuration is an admissible execution fragment of the specification. An implementation is *self-stabilizing* if starting from any configuration, an admissible execution of the implementation eventually reaches a safe configuration. Notice that in the presence of corruptions, if an implementation is self-stabilizing, then any long enough execution fragment of the implementation will eventually have a suffix that looks like the suffix of some correct execution of the specification, until a corruption occurs.

Each of the above services' self-stabilizing implementations will be built on top of self-stabilizing implementations of other services: VtoVComm over the VSA layer, the location service over the VSA layer and VtoVComm service, and EtoEComm over the VSA layer, VtoVComm, and location services. Each self-stabilizing implementation uses lower level services without feedback, so lower level service executions are not influenced by the upper level services. This allows us to guarantee that higher level service implementations are still self-stabilizing through *fair composition* [11].

Our service implementations, starting from an arbitrary system configuration, stabilize within the following times: VtoVComm:  $t_{VtoV} + d$  time after the VSA layer stabilizes ( $t_{VSAcor}$  time), the location service:  $\max(t_{HLS}, 2e + 3t_{VtoV} + t_{hb} + 2d)$  time after VtoVComm stabilizes, and EtoEComm:  $t_{pb} + 2d + 2e + t_{VtoV}$  time after the location service has stabilized.

## 4 VSA to VSA communication (VtoVComm) implementation

The VtoVComm service allows communication of messages between any two VSAs through VtoVsend and VtoVrcv actions, as long as there is a path of non-failed VSAs between them. The VtoVComm service is built on top of the V-bcast service [8], which supports communication between two neighboring VSAs (see Figure 3).

VSA-to-VSA communication is based on a greedy DFS procedure. When a VSA receives a message for which it is not the destination, it chooses a neighboring VSA that is on a shortest path to the destination VSA and forwards the message in a **forward** message to that neighbor. If the VSA does not receive an indication through a **found** message that the message has been delivered to the destination within some bounded amount of time, it then forwards the message to the neighboring VSA on the next shortest path to the destination VSA, and so on. This choice of neighbors is greedy in the sense that the next neighbor chosen to receive the forwarded message is the one on a shortest path to the destination VSA, excluding the neighbors associated with previous tries. The greedy DFS can turn into a flood in pathological situations in which the destination is that last VSA reached.

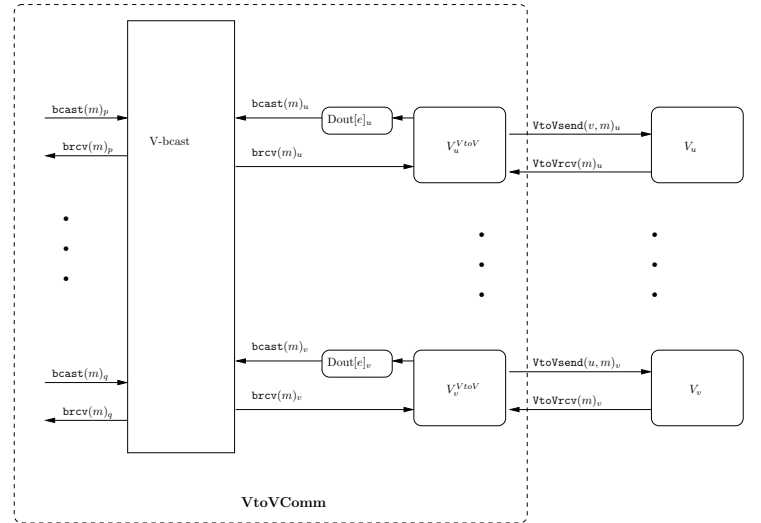


Figure 3: VSA-to-VSA communication (VtoVComm). A VSA at region  $u$  sends a message  $m$  to region  $v$ 's VSA with a  $\text{VtoVsend}(v, m)_u$ . The message is eventually received at region  $v$  by  $\text{VtoVrcv}(m)_v$ .

<p><b>Signature:</b></p> <p>2    <b>Input</b> <math>VtoVsend(d, m)_u, d \in U, m</math> arbitrary</p> <p>4    <b>Input</b> <math>brcv(m)_u, m \in (\{\text{forward}\} \times Msg \times U \times \{u\})</math>  <math>\cup (\{\text{found}\} \times Msg)</math></p> <p>6    <b>Output</b> <math>bcast(m)_u, m</math> arbitrary</p> <p>6    <b>Output</b> <math>VtoVrcv(m)_u, m</math> arbitrary</p> <p>8    <b>Internal</b> <math>DFStimeout(msg)_u, msg \in Msg</math></p> <p>8    <b>Internal</b> <math>DFSclean(msg)_u, msg \in Msg</math></p> <p>10    <math>Msg = M \times U \times U \times \mathbb{R}</math>, of the form <math>\langle m, v2vs, v2vd, ts \rangle</math></p> <p><b>State:</b></p> <p>12    <b>analog</b> <math>now \in \mathbb{R}</math>, the current real time</p> <p>14    <math>bcastq, VtoVrcvq</math>, queues of messages, initially <math>\emptyset</math></p> <p>14    <math>DFStable</math>, a table indexed on message tuples in</p> <p>16    <math>Msg</math> with entries in <math>(nbrs(u) \times 2^{nbrs(u)} \times \mathbb{R})</math>,</p> <p>16    of the form <math>\langle isrc, NbrSet, nbrTO \rangle</math>, initially <math>\emptyset</math></p> <p>18    <math>curNbr \in U</math>, initially <math>\perp</math></p> <p><b>Trajectories:</b></p> <p>20    <b>satisfies</b></p> <p>22    <math>d(now) = 1</math></p> <p>22    <b>constant</b> <math>bcastq, VtoVrcvq, DFStable, curNbr</math></p> <p>24    <b>stops when</b></p> <p>24    Any precondition is satisfied.</p> <p><b>Actions:</b></p> <p>26    <b>Output</b> <math>bcast(m)_u</math></p> <p>28    <b>Precondition:</b></p> <p>28    <math>m \in bcastq</math></p> <p>30    <b>Effect:</b></p> <p>32    <math>bcastq \leftarrow bcastq \setminus \{m\}</math></p> <p>34    <b>Input</b> <math>VtoVsend(d, m)_u</math></p> <p>34    <b>Effect:</b></p> <p>36    <b>if</b> <math>u = d</math> <b>then</b></p> <p>36    <math>VtoVrcvq \leftarrow VtoVrcvq \cup \{m\}</math></p> <p>36    <b>else</b> <math>DFStable(\langle m, u, d, now \rangle) \leftarrow \langle u, nbrs(u), now \rangle</math></p>	<p><b>Internal</b> <math>DFStimeout(msg)_u</math></p> <p><b>Precondition:</b></p> <p>40    <math>DFStable(msg).nbrTO \leq now</math></p> <p>42    <math>\vee DFStable(msg).nbrTO &gt; now + \delta(u, msg.v2vd)</math></p> <p><b>Effect:</b></p> <p>44    <b>if</b> <math>DFStable(msg).NbrSet \neq \emptyset</math> <b>then</b></p> <p>44    <math>curNbr \leftarrow \text{NxtNbr}(DFStable(msg).NbrSet,</math></p> <p>46    <math>DFStable(msg).isrc, u, msg.v2vd)</math></p> <p>46    <math>DFStable(msg).NbrSet \leftarrow DFStable(msg).NbrSet \setminus \{curNbr\}</math></p> <p>48    <math>bcastq \leftarrow bcastq \cup \{\langle \text{forward}, msg, u, curNbr \rangle\}</math></p> <p>48    <math>DFStable(msg).nbrTO \leftarrow now + \delta(u, msg.v2vd)</math></p> <p>50    <b>else</b> <math>DFStable(msg) \leftarrow \text{null}</math></p> <p><b>Input</b> <math>brcv(\langle \text{forward}, msg, isrc, u \rangle)_u</math></p> <p>52    <b>Effect:</b></p> <p>54    <b>if</b> <math>msg.ts \in [now - ttl_{VtoV}, now]</math> <b>then</b></p> <p>54    <b>if</b> <math>u = msg.v2vd</math> <b>then</b></p> <p>56    <math>bcastq \leftarrow bcastq \cup \{\langle \text{found}, msg \rangle\}</math></p> <p>56    <math>VtoVrcvq \leftarrow VtoVrcvq \cup \{msg.m\}</math></p> <p>58    <b>else if</b> <math>DFStable(msg) = \text{null}</math> <b>then</b></p> <p>58    <math>DFStable(msg) \leftarrow \langle isrc, nbrs(u) \setminus \{isrc\}, now \rangle</math></p> <p>60    <b>Input</b> <math>brcv(\langle \text{found}, msg \rangle)_u</math></p> <p>62    <b>Effect:</b></p> <p>62    <b>if</b> <math>DFStable(msg) \neq \text{null}</math> <b>then</b></p> <p>64    <math>DFStable(msg) \leftarrow \text{null}</math></p> <p>64    <b>if</b> <math>u \neq msg.v2vs</math> <b>then</b></p> <p>66    <math>bcastq \leftarrow bcastq \cup \{\langle \text{found}, msg \rangle\}</math></p> <p>68    <b>Output</b> <math>VtoVrcv(m)_u</math></p> <p>70    <b>Precondition:</b></p> <p>70    <math>m \in VtoVrcvq</math></p> <p>72    <b>Effect:</b></p> <p>72    <math>VtoVrcvq \leftarrow VtoVrcvq \setminus \{m\}</math></p> <p>74    <b>Internal</b> <math>DFSclean(msg)_u</math></p> <p>76    <b>Precondition:</b></p> <p>76    <math>DFStable(msg) \neq \text{null} \wedge msg.ts \notin [now - ttl_{VtoV}, now]</math></p> <p>78    <b>Effect:</b></p> <p>78    <math>DFStable(msg) \leftarrow \text{null}</math></p>
---	---

Figure 4: Greedy DFS algorithm at  $V_u^{VtoV}$  for region  $u$ .

Self-stabilization of the algorithm is ensured by the use of a real-time timestamp to identify the version of the DFS. Too old versions are eliminated from the system and new versions are handled as completely new attempts to complete a greedy DFS towards the destination.

We first present a simple greedy DFS algorithm that gradually expands the search until all paths are checked. This algorithm will find a path to the destination if such a path exists throughout the DFS execution. We also present a modification of the algorithm to produce a persistent version of the greedy DFS algorithm in which each VSA repeatedly tries to forward messages along previously unsuccessful VSA paths to take advantage of (possibly temporary) recoveries of VSAs that may result in a viable path [13]. Again, the persistent greedy DFS can turn into a persistent flood in pathological situations in which the destination is the last VSA reached.

#### 4.1 Detailed code description

The following code description refers to the code for VSA  $V_u^{VtoV}$  in Figure 4. The main state variable  $DFStable$  keeps track of information for messages that are still waiting to be delivered. For each such unique message, the table stores the intermediate source  $isrc$  of the message, the set of VSA neighbors  $NbrSet$  of neighbors that have yet to have the message forwarded to them, and a timeout  $nbrTO$  for the neighbor currently being tried for forwarding the message.

A source VSA  $V_u^{VtoV}$  sends a message  $m$  to a destination VSA in region  $d$  using  $VtoVsend(d, m)_u$  (line

```

2 Internal DFStimeout( $msg$ ) $u$ 
3 Precondition:
4    $DFStable(msg).nbrTO \leq now \vee DFStable(msg).nbrTO > now + \delta(u, msg.v2vd)$ 
5 Effect:
6   if  $DFStable(msg).NbrSet \neq \emptyset$  then
7      $curNbr \leftarrow \text{NxtNbr}(DFStable(msg).NbrSet, DFStable(msg).isrc, u, msg.v2vd)$ 
8      $DFStable(msg).NbrSet \leftarrow DFStable(msg).NbrSet \setminus \{curNbr\}$ 
9     for each  $n \in nbr(u) \setminus DFStable(msg).NbrSet$ 
10       $bcstq \leftarrow bcstq \cup \{\{forward, msg, u, n\}\}$ 
11       $DFStable(msg).nbrTO \leftarrow now + \delta(u, msg.v2vd)$ 
12   else  $DFStable(msg) \leftarrow \text{null}$ 

```

Figure 5: The Persistent Greedy DFS algorithm at  $V_u^{VtoV}$  for region  $u$  is the same as the Greedy DFS algorithm, except that the broadcast of a DFS message to  $curNbr$  in the DFStimeout action is replaced with a broadcast to  $curNbr$  and all previously attempted neighbors.

33). If  $u = d$  then  $V_u^{VtoV}$  received  $m$  through  $VtoVrcv(m)_u$  (lines 35-36). Otherwise the destination VSA is another VSA and  $V_u^{VtoV}$  sets the  $DFStable$  mapping of an augmented version of the message,  $\langle m, u, d, now \rangle$ , to  $\langle u, nbrs(u), now \rangle$ . This enables the start of a new DFS execution to forward the message to its destination (line 37).

Whenever the  $nbrTO$  of a message in  $DFStable$  times out, it triggers the forwarding of the message to the next neighbor in the DFS, if possible. If the message hasn't yet been forwarded to all of the relevant neighbors ( $DFStable(msg).NbrSet$  is not empty), then the next neighbor closest to the destination VSA that has not yet had a message forwarded to it,  $curNbr$ , is selected and the message tuple  $msg$  is then forwarded in a **forward** message to it using the V-bcast service (lines 45-48). The timeout variable  $DFStable(msg).nbrTO$  for this attempt at forwarding is set to  $now + \delta(curNbr, msg.v2vd)$  (line 49). If the message has already been forwarded to all the relevant neighbors, then  $DFStable(msg)$  is set to null, indicating that nothing more can be done.

If a message tuple  $msg$  whose destination is  $V_u^{VtoV}$  is received in a **forward** message from  $isrc$ , then VSA  $V_u^{VtoV}$  broadcasts a **found**,  $msg$  message via the V-bcast service and  $VtoVrcv$ 's the message  $msg.m$ . The **found** message notifies neighbors still participating in the DFS for  $msg$  that it has reached its final destination VSA. No forwarding is required (lines 55-57). Otherwise, if  $msg$  is not destined for  $V_u^{VtoV}$  and  $V_u^{VtoV}$  does not already have an entry in  $DFStable$  for  $msg$ , then the message must be forwarded to its destination.  $DFStable(msg)$  is set to  $\langle isrc, nbrs(u) \setminus \{isrc\}, now \rangle$  (line 59), storing the intermediate source, initializing the set of neighbors that have yet to have the message forwarded to them, and setting  $nbrTO$  to  $now$ . Setting  $nbrTO$  to  $now$  immediately enables the DFStimeout action for  $msg$ , triggering the forwarding of  $msg$  to one of  $V_u^{VtoV}$ 's neighbors.

When a **found** message is received for a message tuple  $msg$  that is mapped by  $DFStable$ , the entry in  $DFStable$  is erased, preventing additional forwarding (line 64). If  $u \neq msg.v2vs$  then VSA  $V_u^{VtoV}$  broadcasts a **found** message via the V-bcast service (lines 65-66), notifying neighbors that are still participating for  $msg$  that it has been delivered. Clearly, if  $u = msg.v2vs$ , then no **found** message is required and no further action needs to be taken.

## 4.2 Correctness

We now prove the correctness of the algorithm. Let the source VSA be  $V_s^{VtoV}$ , the destination VSA be  $V_d^{VtoV}$ , the message sent be  $m$ , and a DFS execution  $exe$  from  $V_s^{VtoV}$  to  $V_d^{VtoV}$  be as defined above. We assume a given function  $\delta : \{U\} \times \{U\} \rightarrow \mathcal{N}$ , where  $\delta(x, y)$  is a bound on the time required for a message to arrive from  $x$  to  $y$ . This bound is based both on the distance between  $x$  and  $y$ , and the quality of the communication links in the network. Since the DFS and the  $\delta$  function are just employed to cut down on unneeded retransmission of messages, any non-negative wait time is sufficient for correctness. However, a wait time dependent on hop count between regions will be the most message-efficient. We argue that if no corruption failures occur and the status (failed or non-failed) of every VSA in  $\mathcal{U}$  doesn't change during  $exe$ , then the following holds:

**Lemma 4.1** *If  $V_s^{VtoV}$  is a non-failed VSA that performs a  $VtoVsend(d, m)$  at time  $t$ , and there exists a path of non-failed VSAs between  $V_s^{VtoV}$  and  $V_d^{VtoV}$  from time  $t$  to time  $t + ttl_{VtoV}$ , then  $V_d^{VtoV}$  performs a*

$\mathbf{VtoVrcv}(m)$  in the interval  $[t, t + ttl_{VtoV}]$ , for  $ttl_{VtoV} \geq [e + d + (\max_{u,v \in U} \delta(u, v) \cdot \max_{u \in U} |nbrs(u)| - 1)] \cdot (|U| - 1)$ .

*Proof sketch:* The proof is by induction on the distance  $n$  between  $s$  and  $V_d^{VtoV}$  on the shortest non-deserted path, where the distance is the number of VSAs along the path, including  $V_d^{VtoV}$ . In the case  $n = 0$ , the message  $m$  is destined for the same VSA. According to line 35, the message is  $\mathbf{VtoVrcv}$ 'ed at the VSA.

Let's assume that the lemma holds for every  $n' < n$ .

Let  $n$  be the VSA-distance between  $V_s^{VtoV}$  and  $V_d^{VtoV}$ . There exists a path of non-failed VSAs between  $V_s^{VtoV}$  and  $V_d^{VtoV}$ . Therefore, there exists a VSA  $V_u^{VtoV}$ , which is a neighbor of  $V_s^{VtoV}$ , such that there exists a path of non-failed VSAs between  $V_u^{VtoV}$  and  $V_d^{VtoV}$ . The distance between  $V_u^{VtoV}$  and  $V_d^{VtoV}$  is  $n - 1$ , hence the induction assumption holds for  $V_u^{VtoV}$  and  $V_d^{VtoV}$ . Therefore, a message sent from  $V_u^{VtoV}$  to  $V_d^{VtoV}$  eventually reaches  $V_d^{VtoV}$ . The same assumption holds for  $V_s^{VtoV}$  and  $V_u^{VtoV}$ , therefore,  $V_d^{VtoV}$  receives the message  $m$  sent from region  $s$ . ■

**Lemma 4.2** *The number of times that a message tuple  $msg$  is re-broadcast is bounded.*

*Proof sketch:* The broadcast of a message tuple stops in either of the following cases:

- A **found** message was received for  $msg$ . According to line 62, if the value of  $DFStable(msg)$  was not already null, it gets set to null, preventing  $V_u^{VtoV}$  from doing anything with subsequent **found** messages. If  $V_u^{VtoV}$  was not the original source of  $msg$ , it retransmits **found** for  $msg$  exactly one time. If a **found** for  $msg$  is received again, it will be ignored. A **forward** message for  $msg$  would need to be received again in order to result in any additional **found** messages for  $msg$  at this VSA. This, however, cannot happen since each VSA participating in the DFS waits before triggering new **forward** messages until **found** messages would have been returned.
- For each VSA neighbor, if VSA  $V_u^{VtoV}$  does not receive a **found** message for  $msg$  it will time out via  $nbrTO$ . Once the set of neighbors to be queried is exhausted, the VSA erases the entry for  $msg$  in  $DFStable$ , preventing any additional forwarding by itself. ■

**Lemma 4.3** *Once corruptions stop and the VSA layer has stabilized, it takes up to  $d + ttl_{VtoV}$  time for  $VtoVComm$  to stabilize.*

*Proof sketch:* Any message in the system that is being forwarded by  $VtoVComm$  will be cleaned out of the system if they are older than  $ttl_{VtoV}$  or newer than the current time. As a result, the longest a “bad” message can be in the system is this time, plus up to an additional  $d$  time where it could have been in transmission before being received by a VSA. ■

## 5 Home Location Service (HLS) implementation

The location service, as described in the last section, allows a client to determine a recent region of another alive client. In our implementation, called the *Home Location Service (HLS)*, we accomplish this using *home locations*. Recall that the home locations of a client node  $p$  are  $f + 1$  regions whose VSAs are occasionally updated with  $p$ 's region. The home locations are calculated with a hash function  $h$ , mapping a client's id to a list of VSA regions, and is known to all VSAs. These home location VSAs can then be queried by other VSAs to determine a recent region of  $p$ .

Figure 6 depicts how the VSA abstraction and  $VtoVComm$  are used in HLS. The HLS implementation consists of two parts: a client-side portion and a VSA-side portion.  $C_p^{HL}$  is a subautomaton of client  $p$  that interacts with VSAs to provide HLS. It is responsible for notifying VSAs in its current and neighboring regions which region it is in. Also,  $C_p^{HL}$  handles each request submitted by input  $HLquery(q)_p$  for  $q$ 's region, by broadcasting the query via V-bcast to VSAs  $V_u^{HL}$  in its current and neighboring regions. It translates responses from the VSAs into  $HLreply$  outputs.

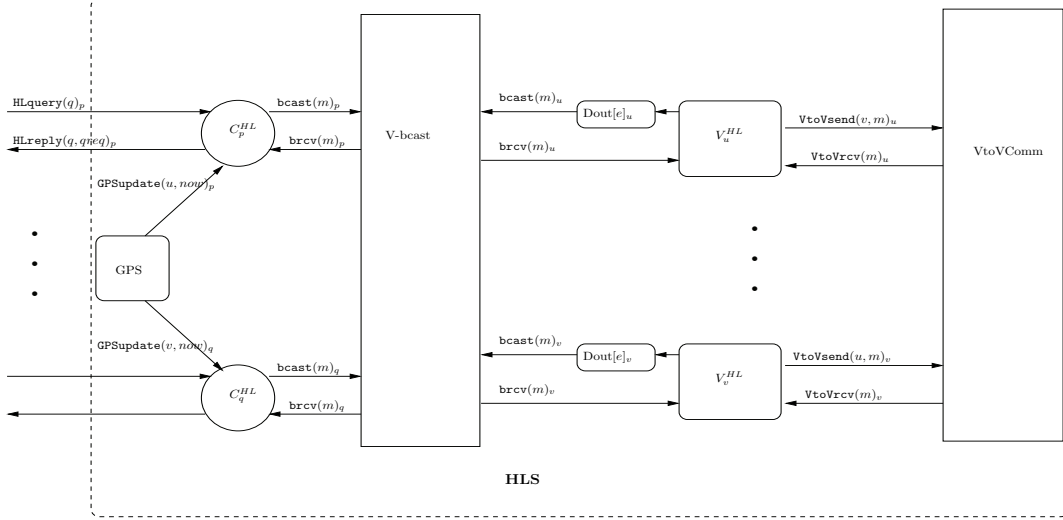


Figure 6: Home Location Service. A client  $p$  can query local VSAs for client  $q$ 's region. The VSAs then query home locations of  $q$ , using  $VtoVComm$ , for a recent region of  $q$ , and return it to  $p$ .

For the VSA-side,  $V_u^{HL}$  and  $V_v^{HL}$  in Figure 6 are home location VSAs corresponding to regions  $u$  and  $v$  of the network; they are subautomata of VSAs  $V_u$  and  $V_v$ .  $V_u^{HL}$  takes a request from a local client for client node  $q$ 's region, calculates  $q$ 's home locations using the hash function, and then sends location queries to the home locations using  $VtoVComm$ . Those virtual automata respond with the region information they have for  $q$ , which is then provided by  $V_u^{HL}$  to the requesting client.  $V_u^{HL}$  also is responsible both for informing the home locations of each client  $p$  located in its region or neighboring regions of  $p$ 's region, and maintaining and answering queries for the regions of clients for which it is a home location.

Time and region information from the GPS oracle is used throughout the HLS algorithm, by clients and VSAs, to timestamp and label information and messages. This information is used to guarantee timeliness of replies from the HLS service, and to stabilize the service after faults. Timestamps are used to determine if information is too old or too new, while region information allows clients and VSAs to know which other clients and VSAs to interact with.

### 5.1 HLS client actions

The code executed by client  $p$ 's  $C_p^{HL}$  is in Figure 7.

Clients receive  $GPSupdates$  every  $\epsilon_{sample}$  time from the GPS automaton (lines 28-33), making them aware of their current region and the time. If a client's region has changed, the client immediately sends a **heartbeat** message with its id, current time and region information. The client periodically reminds its current and neighboring region VSAs of its region by broadcasting additional **heartbeat** messages every  $tll_{hb}$  time, where  $tll_{hb}$  is a known constant (lines 35-39).

$C_p^{HL}$  also handles the  $HLquery(q)$  inputs it receives (line 41). This request for  $q$ 's location is stored in a  $queryq$  table and, once the client knows its own region, translated into a  $\langle clocQuery, q \rangle$  message that is broadcast, together with the VSA region, to local regions' VSAs (lines 45-49). If  $C_p^{HL}$  eventually receives a  $\langle clocReply, q, qreg \rangle$  message from its current or neighboring region's VSA for a client  $q$  in  $queryq$ , indicating that node  $q$  was in region  $qreg$  (lines 51-55), it clears the entry for  $q$  in  $queryq$ , and outputs a  $HLreply(q, qreg)$  of the information (lines 57-61). If the request for  $q$ 's location goes unanswered for more than  $tll_{HLS} - \epsilon_{sample}$  time, then the request has failed and is removed (lines 63-67).

### 5.2 HLS VSA actions

The code for automaton  $V_u^{HL}$  appears in Figure 8.

First, the VSA knows which clients are in its or neighboring regions through **heartbeat** messages. If a VSA hears a **heartbeat** message from a client  $p$  claiming to be in its region or a neighboring region, the

<p><b>Constants:</b></p> <p>2 <math>ttl_{hb}</math></p> <p>4 <math>ttl_{HLS}</math></p> <p><b>Signature:</b></p> <p>6 <b>Input</b> GPSupdate(<math>v, t</math>)<math>_p, v \in U, t \in \mathbb{R}</math></p> <p>8 <b>Input</b> HLquery(<math>q</math>)<math>_p, q \in P</math></p> <p>10 <b>Input</b> brcv(<math>\langle m, u \rangle</math>)<math>_p, m \in (\{\text{clocReply}\} \times P \times U \times U), u \in U</math></p> <p>12 <b>Output</b> bcast(<math>\langle m, reg \rangle</math>)<math>_p, m \in (\text{heartbeat}, now, p) \cup \{\text{clocQuery}\} \times P</math></p> <p>14 <b>Output</b> HLreply(<math>q, v</math>)<math>_p, q \in P, v \in U</math></p> <p>16 <b>Internal</b> queryfail(<math>q</math>)<math>_p, q \in P</math></p> <p><b>State:</b></p> <p>18 <b>analog</b> <math>now \in \mathbb{R}</math>, current real time, initially <math>\perp</math></p> <p>20 <math>hbTO \leq now + ttl_{hb}, \in \mathbb{R}</math>, the next heartbeat time</p> <p>22 <math>reg \in U</math>, the current region, initially <math>\perp</math></p> <p>24 <math>queryq</math>, a table from <math>P</math> to <math>\mathbb{R}</math>, initially <math>\emptyset</math></p> <p>26 <math>queryrcv</math>, a queue of <math>P \times U</math> pairs, initially <math>\emptyset</math></p> <p><b>Trajectories:</b></p> <p>28 <b>satisfies</b></p> <p>30 <math>d(now) = 1</math></p> <p>32 <b>constant</b> <math>hbTO, reg, queryq, queryrcv</math></p> <p><b>stops when</b></p> <p>Any precondition is satisfied.</p> <p><b>Actions:</b></p> <p>34 <b>Input</b> GPSupdate(<math>v, t</math>)<math>_p</math></p> <p>36 <b>Effect:</b></p> <p>38 <math>now \leftarrow t</math></p> <p>40 <b>if</b> <math>reg \neq v</math> <b>then</b></p> <p>42 <math>reg \leftarrow v</math></p> <p>44 <math>hbTO \leftarrow now</math></p>	<p><b>Output</b> bcast(<math>\langle \text{heartbeat}, now, p \rangle, reg \rangle</math>)<math>_p</math></p> <p><b>Precondition:</b></p> <p>36 <math>hbTO \leq now \wedge reg \neq \perp</math></p> <p><b>Effect:</b></p> <p>38 <math>hbTO \leftarrow now + ttl_{hb}</math></p> <p><b>Input</b> HLquery(<math>q</math>)<math>_p</math></p> <p><b>Effect:</b></p> <p>40 <math>queryq(q) \leftarrow \infty</math></p> <p><b>Output</b> bcast(<math>\langle \langle \text{clocQuery}, q \rangle, reg \rangle</math>)<math>_p</math></p> <p><b>Precondition:</b></p> <p>42 <math>reg \neq \perp \wedge queryq(q) &gt; now + ttl_{HLS} - \epsilon_{sample}</math></p> <p><b>Effect:</b></p> <p>44 <math>queryq(q) \leftarrow now + ttl_{HLS} - \epsilon_{sample}</math></p> <p><b>Input</b> brcv(<math>\langle \langle \text{clocReply}, q, qreg \rangle, u \rangle</math>)<math>_p</math></p> <p><b>Effect:</b></p> <p>46 <b>if</b> (<math>u \in nbrs(reg) \cup \{reg\} \wedge queryq(q) \neq null</math>) <b>then</b></p> <p>48 <math>queryrcv \leftarrow queryrcv \cup \{\langle q, qreg \rangle\}</math></p> <p>50 <math>queryq(q) \leftarrow null</math></p> <p><b>Output</b> HLreply(<math>q, qreg</math>)<math>_p</math></p> <p><b>Precondition:</b></p> <p>52 <math>\langle q, qreg \rangle \in queryrcv</math></p> <p><b>Effect:</b></p> <p>54 <math>queryrcv \leftarrow queryrcv \setminus \{\langle q, qreg \rangle\}</math></p> <p><b>Internal</b> queryfail(<math>q</math>)<math>_p</math></p> <p><b>Precondition:</b></p> <p>56 <math>queryq(q) &lt; now</math></p> <p><b>Effect:</b></p> <p>58 <math>queryq(q) \leftarrow null</math></p>
<p>Figure 7: HLS's <math>C_p^{HL}</math> automaton. This client subautomaton serves as a bridge between the client's requests and the VSA layer.</p>	

VSA sends a `locUpdate` message for  $p$ , with  $p$ 's heartbeat timestamp and region, through `VtoVComm` to the VSAs at home locations of client  $p$  (lines 42-46), where home locations are computed using the known hash function  $h$  from  $P \times \{1, \dots, f + 1\}$  to  $U$ .

When a VSA receives one of these `locUpdate` messages for a client  $p$ , it stores both the region indicated in the message as  $p$ 's current region and the attached heartbeat timestamp in its `loc` table (lines 48-51). This location information for  $p$  is refreshed each time the VSA receives a `locUpdate` for client  $p$  with a newer heartbeat timestamp. Since a client sends a `heartbeat` message every  $ttl_{hb}$  time, which can take up to  $d + e$  time to arrive at and trigger a VSA to send a `locUpdate` message through `VtoVComm`, which can take  $ttl_{VtoV}$  time to be delivered at a home location, an entry for client  $p$  is erased if its timestamp is older than  $ttl_{hb} + d + e + ttl_{VtoV}$  (lines 53-57).

The other responsibility of the VSA is to receive and respond to local client requests for location information on other clients. A client  $p$  in a VSA's region or a neighboring region  $v$  can send a query for  $q$ 's current location to the VSA. This is done via a mobile node's broadcast of a  $\langle \langle \text{clocQuery}, q \rangle, v \rangle$  message. When the VSA at region  $u$  receives this query, if no outstanding query for  $q$  exists, it notes the request for  $q$  in  $lquery(q)$ , and sends a `vlocQuery` message to  $q$ 's  $f + 1$  home locations, querying about  $q$ 's location (lines 59-65). Any home location that receives such a message and has an entry for  $q$ 's region responds with a `vlocReply` to the querying VSA with the region (lines 67-70).

If the querying VSA at  $u$  receives a `vlocReply` in response to an outstanding location request for a client  $q$ , it stores the attached region information in  $lquery(q)$  (lines 72-75), broadcasts a `clocReply` message with  $q$  and its region to local clients, and erases the entry for  $lquery(q)$  (lines 77-81). If, however,  $2ttl_{VtoV} + 2e$  time passes since a request for  $q$ 's region was received by a local client and there is no entry for  $q$ 's region,  $lquery(q)$  is just erased (lines 83-87).

<p><b>Constants:</b></p> <p><math>tll_{VtoV}</math></p> <p><math>tll_{hb}</math></p> <p><math>h</math>, a hash function from <math>P \times \{1, \dots, f+1\}</math> to <math>U</math>  such that for <math>p \in P, x, y \in \{1, \dots, f+1\}</math>,  if <math>x \neq y</math>, then <math>h(p, x) \neq h(p, y)</math></p> <p><b>Signature:</b></p> <p><b>Input</b> <math>brcv(\langle m, v \rangle)_u, m \in (\{\text{heartbeat}\} \times \mathbb{R} \times P) \cup (\{\text{clocQuery}\} \times P), v \in U</math></p> <p><b>Input</b> <math>VtoVrcv(\langle v, m \rangle)_u, v \in U, m \in (\{\text{locUpdate}\} \times P \times \mathbb{R}) \cup (\{\text{vlocQuery}\} \times P) \cup (\{\text{vlocReply}\} \times P \times U)</math></p> <p><b>Output</b> <math>bcst(\langle \text{clocReply}, q, qreg, u \rangle)_u, q \in P, qreg \in U</math></p> <p><b>Output</b> <math>VtoVsend(v, m)_u, v \in U</math></p> <p><b>Internal</b> <math>updateHL(q)_u, q \in P</math></p> <p><b>Internal</b> <math>cleanLoc(q)_u, q \in P</math></p> <p><b>Internal</b> <math>cleanLquery(q)_u, q \in P</math></p> <p><b>State:</b></p> <p><math>loc</math>, a table indexed on process ids with entries  from <math>U \times \mathbb{R}^{\geq 0}</math>, of the form <math>\langle reg, ts \rangle</math></p> <p><math>lquery</math>, a table indexed on process ids with entries  from <math>\mathbb{R}^{\geq 0} \times U</math>, of the form <math>\langle to, qreg \rangle</math></p> <p><math>vtovq</math>, a queue of tuples from <math>U \times msg</math>  (Above all initially empty)</p> <p><b>analog</b> <math>now \in \mathbb{R}^{\geq 0}</math>, the current real time</p> <p><b>Trajectories:</b></p> <p><b>satisfies</b></p> <p><math>d(now) = 1</math></p> <p><b>constant</b> <math>loc, lquery, vtovq</math></p> <p><b>stops when</b></p> <p>Any precondition is satisfied.</p> <p><b>Actions:</b></p> <p><b>Output</b> <math>VtoVsend(v, m)_u</math></p> <p><b>Precondition:</b></p> <p><math>\langle v, m \rangle \in vtovq</math></p> <p><b>Effect:</b></p> <p><math>vtovq \leftarrow vtovq \setminus \{\langle v, m \rangle\}</math></p> <p><b>Input</b> <math>brcv(\langle \text{heartbeat}, t, p, v \rangle)_u</math></p> <p><b>Effect:</b></p> <p><b>if</b> <math>(v \in nbrs(u) \cup \{u\} \wedge now - d \leq t \leq now)</math> <b>then</b>  <b>for</b> <math>i = 1</math> <b>to</b> <math>f+1</math>  <math>vtovq \leftarrow vtovq \cup \{\langle h(q, i), \langle v, \text{locUpdate}, q, t \rangle \rangle\}</math></p>	<p><b>Input</b> <math>VtoVrcv(\langle v, \langle \text{locUpdate}, q, t \rangle \rangle)_u</math></p> <p><b>Effect:</b></p> <p><b>if</b> <math>loc(q).ts &lt; t \leq now</math> <b>then</b>  <math>loc(q) \leftarrow \langle v, t \rangle</math></p> <p><b>Internal</b> <math>cleanLoc(q)_u</math></p> <p><b>Precondition:</b></p> <p><math>loc(q).ts \notin [now - tll_{hb} - d - e - tll_{VtoV}, now]</math></p> <p><b>Effect:</b></p> <p><math>loc(q) \leftarrow null</math></p> <p><b>Input</b> <math>brcv(\langle \text{clocQuery}, q, v \rangle)_u</math></p> <p><b>Effect:</b></p> <p><b>if</b> <math>([lquery(q) = null \vee lquery(q).to &lt; now] \wedge v \in nbrs(u) \cup \{u\})</math> <b>then</b>  <math>lquery(q) \leftarrow \langle now + 2tll_{VtoV} + 2e, \perp \rangle</math>  <b>for</b> <math>i = 1</math> <b>to</b> <math>f+1</math>  <math>vtovq \leftarrow vtovq \cup \{\langle h(q, i), \langle u, \langle \text{vlocQuery}, q \rangle \rangle\}</math></p> <p><b>Input</b> <math>VtoVrcv(\langle v, \langle \text{vlocQuery}, q \rangle \rangle)_u</math></p> <p><b>Effect:</b></p> <p><b>if</b> <math>loc(q) \neq null</math> <b>then</b>  <math>vtovq \leftarrow vtovq \cup \{\langle v, \langle u, \langle \text{vlocReply}, q, loc(q).reg \rangle \rangle \rangle\}</math></p> <p><b>Input</b> <math>VtoVrcv(\langle v, \langle \text{vlocReply}, q, qreg \rangle \rangle)_u</math></p> <p><b>Effect:</b></p> <p><b>if</b> <math>lquery(q) \neq null</math> <b>then</b>  <math>lquery(q).qreg \leftarrow qreg</math></p> <p><b>Output</b> <math>bcst(\langle \text{clocReply}, q, lquery(q).qreg, u \rangle)_u</math></p> <p><b>Precondition:</b></p> <p><math>lquery(q).qreg \neq \perp</math></p> <p><b>Effect:</b></p> <p><math>lquery(q) \leftarrow null</math></p> <p><b>Internal</b> <math>cleanLquery(q)_u</math></p> <p><b>Precondition:</b></p> <p><math>lquery(q).to \notin [now, now + 2tll_{VtoV} + 2e]</math></p> <p><b>Effect:</b></p> <p><math>lquery(q) \leftarrow null</math></p>
<p>Figure 8: HLS's <math>V_u^{HL}</math> automaton.</p>	

### 5.3 Correctness

We make the system assumptions described in Section 3. Call  $C_G$  the first global configuration where the system is consistent. For the following two lemmas and theorem, assume we are in a configuration after  $C_G$ , and that no corruption failures occur.

**Lemma 5.1** *For any VSA  $u$ , if there is a request for  $q$ 's region in  $lquery$ , it was submitted through a  $HLquery(q)$  at a client within the last  $\epsilon_{sample} + d + 2tll_{VtoV} + 2e$  time.*

*Proof sketch:* Once a request is submitted by a client to  $C_p^{HL}$ , if the client has not ever received a  $GPSupdate$ , it can take up to  $\epsilon_{sample}$  time for the client to receive one. After the client has received one, it then broadcasts the request to local VSAs, which takes up to  $d$  time to be delivered. VSAs then hold these queries until they expire  $2tll_{VtoV} + 2e$  later. ■

**Lemma 5.2** *Starting  $\epsilon_{sample} + d + e + tll_{VtoV}$  time after client  $p$  enters the system and until  $p$  fails, for each interval of length  $tll_{VtoV} + e$ , all but  $f$  of  $p$ 's home locations will have a non-null  $loc(p)$  entry for the entire interval. If client  $p$  is alive and there is some VSA  $u$  such that  $loc(p)$  is not null,  $p$  was alive and located in  $loc(p).reg$  within the last  $\epsilon_{sample} + d + e + tll_{VtoV}$  time.*

*Proof sketch:* Within  $\epsilon_{sample}$  time of a client entering the system, a  $GPSupdate$  occurs and the client transmits a  $heartbeat$  message. This message can take up to  $d$  time to be received by a nearby VSA, after which it can take  $e + tll_{VtoV}$  time for the VSA to transmit the associated  $locUpdate$  message to the client's home locations and have the message be received, updating any alive home locations'  $loc(p)$  entries. Since for any interval of length  $tll_{hb} + d + 2e + tll_{VtoV}$ , at most  $f$  of the client's home locations can be failed at any point in the interval, all but  $f$  of the client's home locations will receive a  $locUpdate$  message and have a non-null  $loc(p)$  entry, and will remain alive with a non-null  $loc(p)$  entry for at least  $tll_{VtoV} + e$  after the next  $locUpdate$  message is received (within  $tll_{hb} + d + e + tll_{VtoV}$  time after the first was sent). Since this is true for each  $locUpdate$  message, there can only be  $f$  home locations that either do not have a non-null  $loc(p)$  entry or that will not be alive for an additional  $tll_{VtoV} + e$  time.

For the second statement, note that an alive client  $p$  will send a  $heartbeat$  message within  $\epsilon_{sample}$  time of arriving in a region, prompting updates to  $loc(p)$  at alive home locations within  $d + e + tll_{VtoV}$  time. Hence, if a client is alive, any non-null entry for  $loc(p).reg$  can only be as old as  $\epsilon_{sample} + d + e + tll_{VtoV}$ . ■

**Theorem 5.3** *Every client  $p$  searching for a non-failed client  $q$  that has been in the system longer than  $tll_{HLS} + \epsilon_{sample} + d + tll_{VtoV} + e$  time will perform a  $HLreply(q, qreg)$  within time  $tll_{HLS}$ , such that  $q$  was located in region  $qreg$  no more than  $tll_{HLS}$  time ago. No reply will occur if  $q$  has been failed for more than  $tll_{hb} + tll_{HLS} - \epsilon_{sample}$  time. Any reply is in response to a query.*

*Proof sketch:* For the first statement, by the previous lemma, we know that once client  $q$  has been in the system for  $\epsilon_{sample} + d + e + tll_{VtoV}$  time, any queries of its home locations will succeed in producing a result. However, a new  $HLquery$  request “piggybacks” on any prior unexpired  $HLquery$  requests. Since one of these requests could have been initiated just before the client  $q$ 's home locations are updated, we can only guarantee a response will be received for a new request if any outstanding requests will be answered. If the client has been in the system for this total  $tll_{HLS} + d + e + tll_{VtoV}$  time after receiving its first  $GPSupdate$ , then any response to a query can take as much as  $tll_{HLS}$  time:  $\epsilon_{sample}$  time for the querying client to receive its first  $GPSupdate$ ,  $d$  time for the query to be transmitted and received by a local VSA,  $e + tll_{VtoV}$  for the local VSA to query a home location,  $e + tll_{VtoV}$  for the response to arrive at a local VSA,  $e$  time for the local VSA to transmit the response to its requesting clients, and  $d$  time for the transmission to be received and translated into  $HLreplies$  at clients. This total is  $tll_{HLS}$ . As for the age of the response, by the prior lemma, we know that information can only be out of date by  $\epsilon_{sample} + tll_{VtoV} + e + d$  time when a home location responds to a query by another VSA. The response can take  $e + tll_{VtoV}$  time to arrive at the querying VSA, followed by  $e + d$  time for the querying VSA to get the information to the clients that prompted the query. The oldest the information could be is the total.

For the second statement, note that a failed client will not send a  $heartbeat$  message. Since  $loc(p)$  entries are cleared once  $tll_{hb} + d + e + tll_{VtoV}$  time has passed since the  $heartbeat$  message upon which it was based was broadcast, and the information from the entry can only take as much as  $e + tll_{VtoV}$  time to reach a querying VSA and  $e + d$  time to reach any querying clients, the total is the maximum time a  $HLreply$  can occur after the client fails.

For the third statement, note that a query expires after  $tll_{HLS}$  time. Hence, any response generated must be for a query that occurred no more than that time before. ■

**Theorem 5.4** *Starting from an arbitrary configuration, after VtoVComm has stabilized, it takes  $\max(ttl_{HLS}, 2e + 3ttl_{VtoV} + ttl_{hb} + 2d)$  time for HLS to stabilize.*

*Proof sketch:* Once lower levels have stabilized, most client state is made locally consistent within  $\epsilon_{sample}$  time, the time for the client to get a `GPSupdate`. This action resets most variables if the region is updated. The remaining portions of client state are made consistent instantaneously with local correction actions, with the exception of the heartbeat timer and `queryq` variables. The heartbeat timer can only affect operations for at most  $ttl_{hb}$  time. The `queryq` variable can only affect operations for  $ttl_{HLS}$  time, when it would be deleted.

For VSAs, there are two variables that are not instantaneously corrected: `loc` and `lquery`.

The `loc` variable will be consistent within time  $e + 2ttl_{VtoV} + ttl_{hb} + d$ . At worst, there could be a corrupted message that arrives at a VSA after  $ttl_{VtoV}$  time, adding a bad entry that takes  $e + ttl_{VtoV} + ttl_{hb} + d$  time to expire. If the client referred to is in the system, it might not be until the next update after the timestamp of the corrupted message (which could have been delivered as late as  $ttl_{VtoV}$  after corruptions stopped) arrives for the information to be cleaned up. This time is exactly what the offset term for `loc` timeouts describes. Hence, the variable might not be cleaned until  $ttl_{VtoV}$  plus that offset term.

However, there may be responses based on this bad `loc` table information that were sent right at  $e + 2ttl_{VtoV} + ttl_{hb} + d$ , and that take  $e + ttl_{VtoV}$  to arrive at the VSA. The resulting transmission (taking  $d$  time to complete) to local clients is then incorrect. However, those incorrect transmissions cease after the total time  $2e + 3ttl_{VtoV} + ttl_{hb} + 2d$  elapses.

The `lquery` variable is cleaned up within  $ttl_{HLS}$  time. An entry in `lquery` only has a total of  $2ttl_{VtoV} + 2e$  time in the data structure. It could be the case that a spurious request was transmitted in the beginning, which adds  $d$  time. If a region response is received it results in immediate correction of the state through erasure. Hence, the time required to be consistent is the time that it takes for a query to be accounted for.

The maximum of  $ttl_{HLS}$  and  $2e + 3ttl_{VtoV} + ttl_{hb} + 2d$  is the maximum stabilization time. ■

## 5.4 Extensions

Here we briefly describe some possible extensions to our HLS algorithm:

**Home location voting mechanisms:** In systems where corruption failures are limited in number at the VSA level, our implementation could be extended to use a voting mechanism, allowing the “weed-out” of information from corrupted home locations. Rather than querying VSAs waiting for a single region response from a home location VSA, they could wait until the same region is returned from a majority of home locations VSAs. If corruption is limited to some small number of VSAs at a time, but can happen often, then this voting mechanism can be used to provide a stronger location service, immune to these limited number of faults.

**Randomized asymmetric quorums:** It is possible to have asymmetric updates and queries, such as with local updates to close-by VSAs and uniformly selected VSAs or vice versa (the expected number of VSAs that are required to be updated and queried is small, as proved in [22]). Instead of using a predefined set to query, one might use a randomized scheme based on [22], where a random set of regions is chosen for updating and inquiring about the location of a client node. Moreover, we could enhance the scheme in [22] by using a predefined set for location updates (such as the close-by regions) and random set for location queries (or vice versa).

**Attribute queries:** There are scenarios in which one would like to query for client nodes with certain attributes in a geographic area (e.g., a search for a medical doctor that is currently near by). Our scheme supports such queries in a natural way: Attributes can hash to home locations that store tables of clients with the attribute, and their locations. Clients searching for another nearby client with some attribute could then have a local VSA query home locations for the attribute, and select a nearby client from the list that is returned.

## 6 Client end-to-end routing (EtoEComm) implementation

Our implementation of the end-to-end routing service, EtoEComm, uses the location service to discover a recent region location of a destination client node and then uses this location in conjunction with VtoVComm

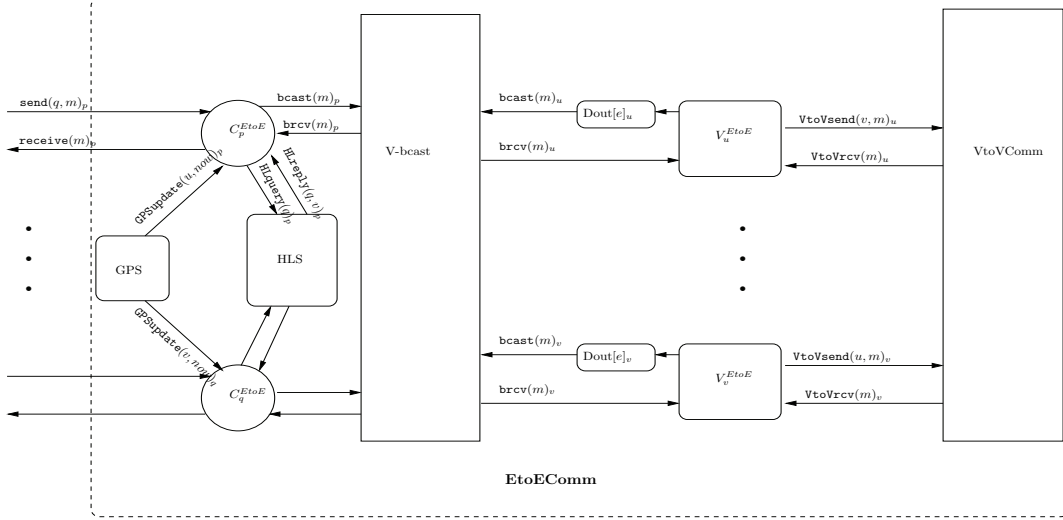


Figure 9: End-to-end routing. A client  $C_p^{E2E}$  can send a message to another client  $C_q^{E2E}$  by querying HLS for  $q$ 's region, and then having local VSAs forward the message to  $q$ 's local regions through VtoVComm. The message is received by those VSAs and broadcast for delivery by  $C_q^{E2E}$ .

to deliver messages (see Figure 9). As in the implementation of the Home Location Service, there are two parts to the end-to-end routing implementation: the client-side portion and the VSA-side portion. Also as in HLS, time and region information from the GPS oracle is used throughout this implementation to timestamp and label information.

The client-side portion  $C_p^{E2E}$  takes a request to send a message to another client  $q$ , queries the HLS for  $q$ 's location, and submits the message to have it sent by a VSA in its current or neighboring regions to  $q$ 's location. It also takes messages originating at other clients and transmitted to it by its current or neighboring regions' VSAs, and delivers them.

The VSA  $V_u^{E2E}$  portion is very simple. A client may send it information to be transmitted to other VSAs, which it forwards through VtoVComm, or another VSA may send it information to be delivered at a client in its own or a neighboring region, which it forwards through V-bcast.

### 6.1 EtoEComm client actions

The signature, state, and actions of  $C_p^{E2E}$  are in Figure 10. The main variable *phbook* is a table, indexed on destination pid, with entries of the form  $\langle reg, ttl, msg \rangle$ . For a client  $q$ , *phbook*( $q$ ).*reg* stores the current region of  $q$  (unless it is unknown, in which case it is  $\perp$ ). The field *ttl* stores a timeout for *phbook*( $q$ ).*reg* if the region of  $q$  is known and stores a timeout for querying for the region if not. The set *msg* stores messages being sent to  $q$ .

The *GPSupdate*( $v, t$ ) action (line 36) results in an update of the client's *reg* variable to the region  $v$  indicated in the action and a reset of the local clock.

A message  $m$  is sent to another client  $q$  via *send*( $q, m$ ) <sub>$p$</sub> . This input to  $C_p^{E2E}$  results in the forwarding of the message to  $p$ 's current region  $u$ 's VSA through *bcast*( $\langle \langle sdata, m, q, phbook(q).reg \rangle, p, u \rangle$ ) if a region *phbook*( $q$ ).*reg* for  $q$  is known (line 44-45), or the saving of the message in *phbook*( $q$ ).*msg*, if the client does not have the location of  $q$  (lines 46-48).

If a recent region for  $q$  is not known,  $C_p^{E2E}$  attempts to discover one. It queries HLS to determine where  $q$  was through the *HLquery*( $q$ ) <sub>$p$</sub>  action (line 50). A timeout for response to the location request, *phbook*( $q$ ).*ttl*, is set for *ttl*<sub>HLS</sub> later. If the timeout expires but no messages are waiting to be sent, *cleanPhbook*( $q$ ) erases the entry, preventing unnecessary *HLquerying* (line 63).

Once a response to an *HLquery*( $q$ ) is received from HLS in the form of *HLreply*( $q, qreg$ ) <sub>$p$</sub>  (line 57), indicating  $q$  was in region  $qreg$ , entry *phbook*( $q$ ).*reg* is updated to  $qreg$  and *phbook*( $q$ ).*ttl* is updated to *now* + *ttl*<sub>pb</sub>,

<p>2 <b>Constants:</b>  <math>ttl_{HLS}</math>  <math>ttl_{pb}</math></p> <p>4 <b>Signature:</b>  6 <b>Input</b> HLreply(<math>q, v</math>)<math>_p, q \in P, v \in U</math>  <b>Input</b> send(<math>q, m</math>)<math>_p, q \in P</math>  8 <b>Input</b> GPSupdate(<math>v, t</math>)<math>_p, v \in U, t \in \mathbb{R}</math>  <b>Input</b> brcv(<math>\langle \langle rdata, m \rangle, p, u \rangle</math>)<math>_p, u \in U</math>  10 <b>Output</b> bcast(<math>m</math>)<math>_p</math>  <b>Output</b> HLquery(<math>q</math>)<math>_p, q \in P</math>  12 <b>Output</b> receive(<math>m</math>)<math>_p</math>  <b>Internal</b> cleanPhbook(<math>q</math>)<math>_p, q \in P</math></p> <p>14 <b>State:</b>  16 <b>analog</b> <math>now \in \mathbb{R}</math>, current real time, initially <math>\perp</math>  <math>reg \in U</math>, the current region, initially <math>\perp</math>  18 <math>phbook</math>, a table indexed on process id with entries from  <math>U \times \mathbb{R} \times 2^{msg}</math>, of the form <math>\langle reg, ttl, msg \rangle</math>, initially <math>\emptyset</math>  20 <math>sdataq, deliverq</math>, queues of messages, initially <math>\emptyset</math></p> <p>22 <b>Trajectories:</b>  <b>satisfies</b>  24 <math>d(now) = 1</math>  <b>constant</b> <math>reg, phbook, sdataq, deliverq</math>  26 <b>stops when</b>  Any precondition is satisfied.</p> <p>28 <b>Actions:</b>  30 <b>Output</b> bcast(<math>\langle \langle sdata, m, q, qreg \rangle, p, reg \rangle</math>)<math>_p</math>  <b>Precondition:</b>  32 <math>\langle m, q, qreg \rangle \in sdataq \wedge reg \neq \perp</math>  <b>Effect:</b>  34 <math>sdataq \leftarrow sdataq \setminus \{ \langle m, q, qreg \rangle \}</math></p> <p>36 <b>Input</b> GPSupdate(<math>v, t</math>)<math>_p</math>  <b>Effect:</b>  38 <math>now \leftarrow t</math>  <b>if</b> <math>reg \neq v</math> <b>then</b>  40 <math>reg \leftarrow v</math></p>	<p><b>Input</b> send(<math>q, m</math>)<math>_p</math>  <b>Effect:</b>  <b>if</b> (<math>phbook(q).reg \neq \perp \wedge phbook(q).ttl \geq now</math>) <b>then</b>  44 <math>sdataq \leftarrow sdataq \cup \{ \langle m, q, phbook(q).reg \rangle \}</math>  <b>else if</b> (<math>phbook(q) = null \vee phbook(q).ttl &lt; now</math>) <b>then</b>  46 <math>phbook(q) \leftarrow \langle \perp, \perp, \{m\} \rangle</math>  <b>else</b> <math>phbook(q).msg \leftarrow phbook(q).msg \cup \{m\}</math>  48</p> <p><b>Output</b> HLquery(<math>q</math>)<math>_p</math>  <b>Precondition:</b>  50 <math>phbook(q) = \langle \perp, ttl, m \neq \emptyset \rangle</math>  <math>\wedge (ttl = \perp \vee ttl &gt; now + ttl_{HLS})</math>  <b>Effect:</b>  52 <math>phbook(q).ttl \leftarrow now + ttl_{HLS}</math>  54</p> <p><b>Input</b> HLreply(<math>q, qreg</math>)<math>_p</math>  <b>Effect:</b>  56 <b>for each</b> <math>m \in phbook(q).msg</math>  <math>sdataq \leftarrow sdataq \cup \{ \langle m, q, qreg \rangle \}</math>  60 <math>phbook(q) \leftarrow \langle qreg, now + ttl_{pb}, \emptyset \rangle</math>  62</p> <p><b>Internal</b> cleanPhbook(<math>q</math>)<math>_p</math>  <b>Precondition:</b>  64 <math>phbook(q) = \langle qreg, ttl, msg \rangle \wedge [ (qreg = \perp \wedge msg = \emptyset) \vee (qreg \neq \perp \wedge [ttl &gt; now + ttl_{pb} \vee msg \neq \emptyset]) \vee ttl &lt; now ]</math>  66 <b>Effect:</b>  68 <math>phbook(q) \leftarrow null</math></p> <p><b>Input</b> brcv(<math>\langle \langle rdata, m \rangle, p, u \rangle</math>)<math>_p</math>  <b>Effect:</b>  70 <b>if</b> <math>u \in \{reg\} \cup nbrs(reg)</math>  <math>deliverq \leftarrow deliverq \cup \{m\}</math>  72</p> <p><b>Output</b> receive(<math>m</math>)<math>_p</math>  <b>Precondition:</b>  74 <math>m \in deliverq</math>  <b>Effect:</b>  76 <math>deliverq \leftarrow deliverq \setminus \{m\}</math>  78</p>
--	--

Figure 10: EtoEComm's  $C_p^{E2E}$  automaton.

storing the location of  $q$  and setting a timeout for use of the location information. For each message waiting to be sent to  $q$  in queue  $phbook(q).msg$ , the message, with the location information for the destination, is forwarded to  $p$ 's current and neighboring regions' VSAs through a  $bcast(\langle \langle sdata, m, q, qreg \rangle, p, u \rangle)$  (lines 59-60, 30-34).

Messages for client  $p$  from other clients are received from  $p$ 's current region or a neighboring region  $v$ 's VSA through  $brcv(\langle \langle rdata, m \rangle, p, v \rangle)_p$  (line 70). The message  $m$  is subsequently delivered through the output  $receive(m)_p$  (line 75).

## 6.2 EtoEComm VSA actions

The signature, state, and actions of  $V_u^{E2E}$  are in Figure 11.

The receipt of a message  $m$  to be sent from a client  $p$  to  $q$  at  $qreg$  through  $brcv(\langle \langle sdata, m, q, qreg \rangle, p, v \rangle)$ ,  $v$  either  $u$  or a neighbor (line 33) results in the subsequent forwarding of the message to the virtual automata at regions in  $calcregs(qreg)$  and their neighboring regions, via the virtual automata communication action  $VtoVsend(qreg, \langle data, m, q \rangle)_u$  (line 33-38). The set  $calcregs(qreg)$  contains the regions that  $q$  could occupy by the time the message is delivered to it (since we do not require the client to be stationary during execution of the algorithm). As will be seen shortly, the definition of  $calcregs$  is dependent on assumptions about client mobility.

Likewise, the receipt, via  $VtoVrcv(\langle data, m, p \rangle)_u$  (line 40), of message  $m$  intended for client  $p$  results in

<p><b>Signature:</b></p> <p>2    <b>Input</b> <math>VtoVrcv(\langle \text{data}, m, p \rangle)_u, p \in P</math></p> <p>3    <b>Input</b> <math>brcv(\langle \langle \text{sdata}, m, q, qreg \rangle, p, v \rangle)_u, p, q \in P, qreg, v \in U</math></p> <p>4    <b>Output</b> <math>bcast(m)_u</math></p> <p>5    <b>Output</b> <math>VtoVsend(v, m)_u, v \in U</math></p> <p>6</p> <p><b>State:</b></p> <p>8    <math>vtovq</math>, a queue of tuples from <math>U \times msg</math>, initially <math>\emptyset</math></p> <p>9    <math>bcastq</math>, a queue of messages, initially <math>\emptyset</math></p> <p>10</p> <p><b>Trajectories:</b></p> <p>12    <b>satisfies</b></p> <p>13       <b>constant</b> <math>vtovq, bcastq</math></p> <p>14    <b>stops when</b></p> <p>15       Any precondition is satisfied.</p> <p>16</p> <p><b>function</b> <math>calcregs(v: U): 2^U =</math></p> <p>18    <b>return</b> <math>nbrs(v) \cup \{v\}</math></p>	<p><b>Actions:</b></p> <p><b>Output</b> <math>bcast(m)_u</math></p> <p><b>Precondition:</b></p> <p>22    <math>m \in bcastq</math></p> <p><b>Effect:</b></p> <p>24    <math>bcastq \leftarrow bcastq \setminus \{m\}</math></p> <p>26</p> <p><b>Output</b> <math>VtoVsend(v, m)_u</math></p> <p><b>Precondition:</b></p> <p>28    <math>\langle v, m \rangle \in vtovq</math></p> <p><b>Effect:</b></p> <p>30    <math>vtovq \leftarrow vtovq \setminus \{\langle qreg, m \rangle\}</math></p> <p>32</p> <p><b>Input</b> <math>brcv(\langle \langle \text{sdata}, m, q, qreg \rangle, p, v \rangle)_u</math></p> <p><b>Effect:</b></p> <p>34    <b>if</b> <math>v \in nbrs(u) \cup \{u\}</math> <b>then</b></p> <p>35       <b>let</b> <math>regions = calcregs(qreg)</math> <b>in</b></p> <p>36          <b>for each</b> <math>v \in regions \cup nbrs(regions)</math></p> <p>37             <math>vtovq \leftarrow vtovq \cup \{\langle qreg, \langle \text{data}, m, q \rangle \rangle\}</math></p> <p>38</p> <p><b>Input</b> <math>VtoVrcv(\langle \text{data}, m, p \rangle)_u</math></p> <p><b>Effect:</b></p> <p>40    <math>bcastq \leftarrow bcastq \cup \{\langle \langle \text{rdata}, m \rangle, p, u \rangle\}</math></p> <p>42</p>
---	--

Figure 11: EtoEComm's  $V_u^{E2E}$  automaton.

the forwarding of the message to  $p$  via  $bcast(\langle \langle \text{rdata}, m \rangle, p, u \rangle)_u$  (line 42).

### 6.3 Correctness

We make the system assumptions described in Section 3. Correctness of the EtoEComm implementation is dependent on assumptions about client mobility and the definition of the function `calcregs`, used in the EtoEComm VSA algorithm. We can prove correctness under either of the following two conditions:

- (1) `calcregs(qreg)` returns the set containing  $qreg$  and its neighbors, and each client remains in a region at least  $\epsilon_{sample} + 3ttl_{VtoV} + 5e + 4d + ttl_{pb}$  time before moving to a neighboring region, or
- (2) `calcregs(qreg)` returns the set containing  $qreg$  and each region  $v$  such that the supremum distance between any two points in  $v$  and  $qreg$  is at most  $v_{max} \cdot (\epsilon_{sample} + 3ttl_{VtoV} + 5e + 4d + ttl_{pb})$ .

We then outline correctness for EtoEComm under these assumptions. For the first lemma and theorem, assume we start in a safe configuration and no corruption failures occur.

**Lemma 6.1** *Consider an alive client  $q$  such that some other client  $p$  has a non-null, non- $\perp$  entry for  $phbook(q).reg$ . If  $q$  does not fail for an additional  $2d + 2e + ttl_{VtoV}$  time, then at any point in that interval,  $q$  will be located in a region in  $calcregs(phbook(q).reg)$ .*

*Proof sketch:* First, we note that a non-null, non- $\perp$  entry  $phbook(q).reg$  has information that is at most  $\epsilon_{sample} + 2ttl_{VtoV} + 3e + 2d$  out-of-date (from HLS) when it is first installed, after which it is saved for an additional  $ttl_{pb}$  time.

If we are assuming condition 1, client  $q$  must be in the region indicated, or a neighboring region, and will remain in those regions for an additional  $2d + 2e + ttl_{VtoV}$  time. If we are assuming condition 2, at any point up to  $2d + 2e + ttl_{VtoV}$  later, client  $q$  can be in any region reachable from  $qreg$  in the total  $\epsilon_{sample} + 3ttl_{VtoV} + 5e + 4d + ttl_{pb}$  time, when traveling at speed  $v_{max}$ . ■

**Theorem 6.2** *Consider a client  $p$  that performs a  $send(q, m)$ , and does not change regions for  $ttl_{HLS}$  time. If client  $q$  has been in the system for  $ttl_{HLS} + \epsilon_{sample} + d + ttl_{VtoV} + e$  time and does not fail, then  $q$  will perform a  $receive(m)$  within  $ttl_{HLS} + 2d + 2e + ttl_{VtoV}$  time. If a client receives a message, it must previously have been sent to it.*

**Theorem 6.3** *Starting from an arbitrary configuration, after HLS has stabilized, it takes  $ttl_{pb} + 2d + 2e + ttl_{VtoV}$  time for EtoEComm to stabilize.*

*Proof sketch:* Bad region information can be in  $phbook$  for up to  $ttl_{pb}$  time, and messages sent using this information are not delivered and cleared until up to  $d + e + ttl_{VtoV} + e + d$  later. At the same time, while HLS has been stabilizing,  $phbook$ 's message collection can take up to  $ttl_{HLS}$  time to be cleared. The maximum of these quantities is the time for EtoEComm to stabilize. ■

## 6.4 Extensions

Here we briefly describe some possible extensions to our EtoEComm algorithm:

**Routing optimizations:** Once the location of a client is known, communication with the client can be continued directly, and movements during the conversation may be piggy-backed on the information transferred in order to update the destination according to the move (as suggested [12]). We also note that we can use an embedded tree location scheme such as the one in [12], implemented by virtual automata, where intermediate tree nodes are also mapped to regions.

**Sleeping client messaging service:** Mobile clients might be able to shut down to conserve power. We could guarantee that a sleeping client eventually receives messages intended for it by having local VSAs save the messages. The VSAs then, at predefined times, broadcast the messages. Sleeping clients awake for these broadcasts, receive their messages, and can go to sleep again afterwards.

## 7 Concluding remarks

We described how both the GPS oracle and the VSA programming layer could help implement self-stabilizing geocast routing, location management, and end-to-end routing services. The self-stabilizing VSA layer provides a virtual fixed infrastructure useful for solving a variety of problems. It acts as a fault-tolerant, self-stabilizing building block for services, allowing applications to be built for mobile networks as though base stations existed for mobile clients to interact with.

The GPS oracle's frequently refreshed and reliable timing and location information made providing self-stabilization easier. We believe the paradigm of an external service providing reliable information that can be used in a self-stabilizing service implementation is an especially important and relevant one in mobile networks. Mobile networks demonstrate many properties that naturally require self-stabilizing implementations, such as a need for self-configuration, or the possibility of unpredictable kinds of failures, but also often have access to reliable external knowledge that can act as a source of shared consistency in the network; here, accurate region knowledge allowed nodes to determine who they should be communicating with (current region and neighboring region nodes), and time information allowed them to order messages and assess timeliness of information.

## References

- [1] Abraham, I., Dolev, D., and Malkhi, D., "LLS: A Locality Aware Location Service for Mobile Ad Hoc Networks", *Proceedings of the DIALM-POMC Joint Workshop on Foundations of Mobile Computing (DIALM-POMC)*, pp. 75-84, 2004.
- [2] Arora, A., Demirbas, M., Lynch, N., and Nolte, T., "A Hierarchy-based Fault-local Stabilizing Algorithm for Tracking in Sensor Networks", *8th International Conference on Principles of Distributed Systems (OPODIS)*, 2004.
- [3] Camp, T., Liu, Y., "An adaptive mesh-based protocol for geocast routing", *Journal of Parallel and Distributed Computing: Special Issue on Mobile Ad-hoc Networking and Computing*, pp. 196-213, 2002.
- [4] Dijkstra, E.W., "Self stabilizing systems in spite of distributed control", *Communications of the ACM*, pp. 643-644, 1974.
- [5] Dolev, S., *Self-Stabilization*, MIT Press, 2000.
- [6] Dolev, S., Gilbert, S., Lynch, N., Schiller, E., Shvartsman, A., and Welch, J., "Virtual Mobile Nodes for Mobile Ad Hoc Networks", *International Conference on Principles of Distributed Computing (DISC)*, pp. 230-244, 2004.
- [7] Dolev, S., Gilbert, S., Lynch, N., Shvartsman, A., Welch, J., "GeoQuorums: Implementing Atomic Memory in Ad Hoc Networks", *17th International Conference on Principles of Distributed Computing (DISC)*, Springer-Verlag LNCS:2848, pp. 306-320, 2003. Also to appear in *Distributed Computing*.
- [8] Dolev, S., Gilbert, S., Lahiani, L., Lynch, N., and Nolte, T., "Timed Virtual Stationary Automata for Mobile Networks", Technical Report MIT-LCS-TR-979a, MIT CSAIL, Cambridge, MA 02139, 2005.
- [9] Dolev, S., Gilbert, S., Lahiani, L., Lynch, N., and Nolte, T., "Brief Announcement: Virtual Stationary Automata for Mobile Networks", *Proceedings of the 24th Annual ACM Symposium on Principles of Distributed Computing (PODC)*, pp. 323, 2005.
- [10] Dolev, S., Herman, T., and Lahiani, L., "Polygonal Broadcast, Secret Maturity and the Firing Sensors", *Third International Conference on Fun with Algorithms (FUN)*, pp. 41-52, May 2004. Also to appear in *Ad Hoc Networks Journal*, Elsevier.
- [11] Dolev, S., Israeli, A., and Moran, S., "Self-Stabilization of Dynamic Systems Assuming only Read/Write Atomicity", *Proceeding of the ACM Symposium on the Principles of Distributed Computing (PODC 90)*, pp. 103-117. Also in *Distributed Computing* 7(1): 3-16 (1993).

- [12] Dolev, S., Pradhan, D.K., and Welch, J.L., "Modified Tree Structure for Location Management in Mobile Environments", *Computer Communications*, Special issue on mobile computing, Vol. 19, No. 4, pp. 335-345, April 1996. Also INFOCOM 1995, Vol. 2, pp. 530-537, 1995.
- [13] Dolev, S. and Welch, J.L., "Crash Resilient Communication in Dynamic Networks", *IEEE Transactions on Computers*, Vol. 46, No. 1, pp.14-26, January 1997.
- [14] Haas, Z.J. and Liang, B., "Ad Hoc Mobility Management With Uniform Quorum Systems", *IEEE/ACM Trans. on Networking*, Vol. 7, No. 2, pp. 228-240, April 1999.
- [15] Hubaux, J.P., Le Boudec, J.Y., Giordano, S., and Hamdi, M., "The Terminodes Project: Towards Mobile Ad-Hoc WAN", *Proceedings of MOMUC*, pp. 124-128, 1999.
- [16] Imielinski, T., and Badrinath, B.R., "Mobile wireless computing: challenges in data management", *Communications of the ACM*, Vol. 37, Issue 10, pp. 18-28, 1994.
- [17] Karp, B. and Kung, H. T., "GPSR: Greedy Perimeter Stateless Routing for Wireless Networks", *Proceedings of the 6th Annual International Conference on Mobile Computing and Networking*, pp. 243-254, SCM Press, 2000.
- [18] Kaynar, D., Lynch, N., Segala, R., and Vaandrager, F., "The Theory of Timed I/O Automata", Technical Report MIT-LCS-TR-917a, MIT LCS, 2004.
- [19] Kuhn, F., Wattenhofer, R., Zhang, Y., Zollinger, A., "Geometric Ad-Hoc Routing: Of Theory and Practice", *Proceedings of the 22nd Annual ACM Symposium on Principles of Distributed Computing (PODC)*, pp. 63-72, 2003.
- [20] Kuhn, F., Wattenhofer, R., and Zollinger, A., "Asymptotically Optimal Geometric Mobile Ad-Hoc routing", *Proceedings of the 6th International Workshop on Discrete Algorithms and Methods for Mobile Computing and Communications (Dial-M)*, pp. 24-33, ACM Press, 2002.
- [21] Li, J., Jannotti, J., De Couto, D.S.J., Karger, D.R., and Morris, R., "A Scalable Location Service for Geographic Ad Hoc Routing", *Proceedings of Mobicom*, pp. 120-130, 2000.
- [22] Malkhi, D., Reiter, M., and Wright, R., "Probabilistic Quorum Systems", *Proceeding of the 16th Annual ACM Symposium on the Principles of Distributed Computing (PODC 97)*, pp. 267-273, Santa Barbara, CA, August 1997.
- [23] Nath, B., Niculescu, D., "Routing on a curve", *ACM SIGCOMM Computer Communication Review*, pp. 155-160, 2003.
- [24] Navas, J.C., Imielinski, T., "Geocast- geographic addressing and routing", *Proceedings of the 3rd MobiCom*, pp. 66-76, 1997.
- [25] Ratnasamy, S., Karp, B., Yin, L., Yu, F., Estrin, D., Govindan, R., and Shenker, S., "GHT: A Geographic Hash Table for Data-Centric Storage", *First ACM International Workshop on Wireless Sensor Networks and Applications*, pp. 78-87, 2002.