



Agent-Based Forwarding Strategies for Reducing Location Management Cost in Mobile Networks

ING-RAY CHEN

*Department of Computer Science, Virginia Polytechnic Institute and State University, Northern Virginia Center,
7054 Haycock Road, Falls Church, VA 22043, USA*

TSONG-MIN CHEN and CHIANG LEE

Computer Science and Information Engineering Department, National Cheng Kung University, Tainan, Taiwan

Abstract. An important issue in location management for dealing with user mobility in wireless networks is to reduce the cost associated with location updates and searches. The former operation occurs when a mobile user moves to a new location registration area and the network is being informed of the mobile user's current location; the latter operation occurs when there is a call for the mobile user and the network must deliver the call to the mobile user. In this paper, we propose and analyze a class of new agent-based forwarding schemes with the objective to reduce the location management cost in mobile wireless networks. We develop analytical models to compare the performance of the proposed schemes with existing location management schemes to demonstrate their feasibility and also to reveal conditions under which our proposed schemes are superior to existing ones. Our proposed schemes are particularly suitable for mobile networks with switches which can cover a large number of location registration areas.

Keywords: location management algorithms, personal communication service networks, two-location algorithm, agent-based forwarding algorithm, Markov chains

1. Introduction

Personal Communication Services (PCS) networks provide wireless communication services to mobile users at any time and any place. In order to search a mobile user (MU) in response to a call, an efficient location management scheme is important. A location management scheme must handle two issues efficiently: location updates and searches. The former operation occurs when a MU moves to a new Visitor Location Register (VLR) area and the network is being informed of where the MU is currently located; the latter operation occurs when there is a call for the MU and the network must route the call to the MU.

Under the basic IS-41 [3] scheme, a MU is permanently registered under a location register called the Home Location Register (HLR). Whenever a MU moves into a new VLR, the MU's HLR is informed of the location change so that it keeps track of the current VLR exactly. When there is a call asking for the MU, the system must query the HLR to get the location of the called MU's current VLR. This scheme is also known as the basic HLR/VLR scheme.

In recent years, several location algorithms have been proposed to reduce the location update cost in mobile wireless networks. When the frequency of the incoming calls is higher than the mobile user's mobility rate, that is, when the call-to-mobility ratio (CMR) is high, the location cache scheme [6] is proposed to reduce the number of searching operations. When CMR is low, on the other hand, it is reported that the forwarding and resetting algorithm (FRA) [2,7], the alternative location strategy (ALS) [13] and the

two-location algorithm (TLA) [9,10] can be used to reduce the location update cost. The CMR value is presumably a per-MU measure since the mobility pattern and the call frequency vary from one MU to another. Therefore, different location algorithms can be adopted depending on the CMR value of the MU user in question. To compare the performance of location management algorithms, one should compare them under the same network environment setting for the same CMR value.

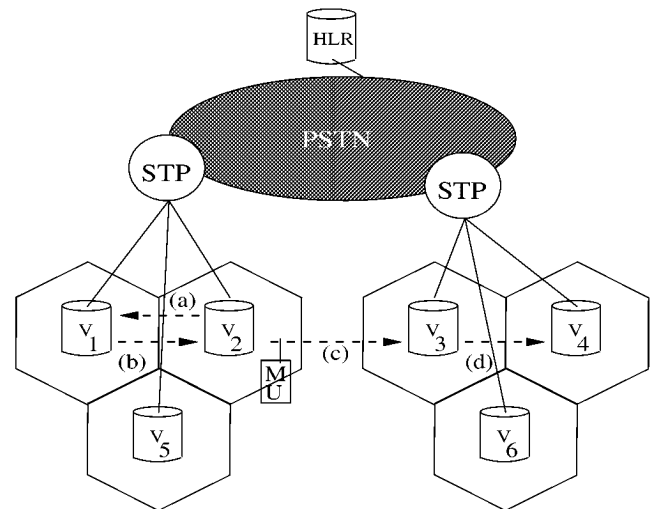
In this paper, we propose and analyze a class of agent-based forwarding strategies that extend and enhance the TLA scheme proposed by Lin [9] by applying the local agent concept [1,4,13] and the forwarding concept [2,7]. Our strategies are intended for MUs with relatively low CMR values, as TLA, ALS and FRA were designed for. The TLA scheme originally proposed by Lin [9] was introduced to reduce the location update cost by recording two most recently visited VLRs in both the HLR database and the MU database. In some cases such as the MU moves back and forth between two VLRs, the location database in the HLR does not have to be updated, thus saving the network cost due to location updates. The saving is especially significant when the CMR value of the MU is low. Of course, a penalty has to be paid when there is a "location miss" since the HLR's database is not up-to-date all the time and two searches instead of just one may be needed to locate the MU. Therefore, if the CMR value is large enough, the cost saving in location updates will eventually be outweighed by the high cost in location searches.

In this paper, we show that the performance of our proposed agent-based forwarding strategies can perform better than the TLA scheme which has been shown to perform very well at low CMR values [9,10]. We present two strategies in this context, namely, the Agent-Based Two-Location Algorithm (ATLA) and the Cross-Update Agent-Based Two-Location Algorithm (CATLA) strategies. The basic idea is to use a VLR as the local agent on behalf of the HLR while applying the TLA scheme to the local agent instead of to the HLR. This effectively reduces the VLR–HLR communication cost to the VLR–VLR communication cost. The local agent of a MU may change as the MU moves across the VLR boundary, depending on the rules being applied. CATLA is different from ATLA by virtue of one rule: the local agent is replaced whenever the MU moves to a new VLR covered by a different network switch from the one that covers the current VLR. Our performance analysis is based on analytical modeling. Specifically, we develop three separate Markov models to describe the behavior of the network under TLA, ATLA, and CATLA, respectively. We then give values to model parameters for the same network setting (by means of a network coverage model) and workload condition, and then assess if our scheme works better. In this paper, we use the combined cost due to location searches and updates between two successive calls to a MU as a metric to compare our schemes with existing ones, as having been done in [10].

The rest of the paper is organized as follows. Section 2 introduces the background. We first describe existing schemes, including the IS-41 basic HLR/VLR and TLA schemes. Then we describe our proposed ATLA and CATLA schemes in detail to deal with mobility-induced location management issues. In section 3, we develop three Markov models to describe the behavior of the network operating under TLA, ATLA and CATLA separately. Section 4 parameterizes these Markov models based on a hexagonal network coverage model so as to compare the performance of ATLA and CATLA schemes with IS-41 and TLA under identical conditions. Section 5 summarizes the paper and outlines some possible future research areas.

2. Background

A Personal Communication Services (PCS) network architecture is shown as in figure 1. The regions covered by the wireless communicable area are divided into many registration areas (RAs) [5]. Assume that each RA has its own VLR. Typically, in a network as such there are switches connecting the HLR to VLRs and also between VLRs. For example, V_1 , V_2 and V_5 are three VLRs under one switch transfer point (STP), while V_3 , V_4 and V_6 are under another STP. STPs are connected by the public switch telephone system (PSTN). Figure 1 also shows a hierarchical PCS network as discussed in [6] in which there is only one HLR for each MU, but the MU may go to different registration areas under different VLRs (marked V_0 , V_1 , V_2 , etc. in figure 1). As we can see, when the HLR and VLRs are in different seg-



PSTN: public switched telephone network

STP : signal transfer point

HLR : home location register

V_i : visitor location register

MU : mobile user

Figure 1. A mobile communication network.

ments of the network, the location update cost can be substantial under the basic HLR/VLR scheme, especially when the MU moves frequently since every move involves a connection cost between the HLR and the new VLR which the MU moves into.

2.1. IS-41 basic HLR/VLR scheme

This basic HLR/VLR scheme exists in IS-41 [3] in the United States and GSM [11] in Europe. Under the basic HLR/VLR scheme, a MU is permanently registered under a location register called the home location register (HLR). When the MU moves to a new location area, it reports to the new VLR which in turn sends a location update message to the HLR that the MU is now under its area. When there is a caller asking for the MU, the caller first checks the local VLR for the MU's profile. If the profile is not in the local VLR, a message is sent to the HLR querying the MU's location. The MU's HLR verifies the current VLR location and then returns the location to the caller.

2.2. Two-location algorithm (TLA)

Under the two-location algorithm (TLA) scheme, a MU and its HLR each keep a *location table* to store two recently visited VLRs. When a MU moves to a new VLR which is not one of the two recently visited VLRs recorded in the table, an update operation is initiated by the mobile user so that both the location tables in the mobile unit and in HLR are updated. Figure 1 shows a MU moving from V_2 to V_1 and then back to V_2 ; after moving back to V_2 , the MU moves to

Table 1
Databases stored in the HLR and MU under TLA.

Movement	HLR location table	MU location table	Update cost
(a) $V_2 \rightarrow V_1$	(V_1, V_2)	(V_1^*, V_2)	VLR–HLR
(b) $V_1 \rightarrow V_2$	(V_1, V_2)	(V_2^*, V_1)	No
(c) $V_2 \rightarrow V_3$	(V_3, V_2)	(V_3^*, V_2)	VLR–HLR
(d) $V_3 \rightarrow V_4$	(V_4, V_3)	(V_4^*, V_3)	VLR–HLR

V_3 and finally to V_4 . In table 1, we illustrate TLA by showing the contents of the databases maintained by the HLR and MU upon movements (a)–(d). We use “*” to mark the current VLR being visited by the MU. Initially, assume that the MU is in V_2 . When the MU makes movement (a), an update operation is performed with the HLR and the location tables of HLR and MU are both updated to (V_1^*, V_2) , with V_1 being the current VLR. For the HLR’s database, this sequence also corresponds to the search sequence. After movement (b), the MU moves back to the previously visited VLR, i.e., V_2 , but the HLR database does not need to be updated since V_2 is already in the HLR’s location table. Note that in this state, the HLR’s database is inconsistent, i.e., it still keeps V_1 as the current VLR of the MU and thus will search V_1 first to service a call. In this case, if a call arrives then there will be a location miss first at V_1 , and a double search cost will incur, although the MU will eventually be found at V_2 . Continuing with our movement example, the MU subsequently makes movements (c) and (d), each of which also requires a location update to the HLR database.

2.3. ATLA

In this section, we describe a proposed new scheme called ATLA. This agent-based forwarding strategy exploits the locality properties of the MU’s movement and access history pattern. The idea is as follows. The HLR at all times only points to a single (but changeable) VLR as in the IS-41 basic HLR/VLR scheme. We will call this VLR the local agent of the MU. As in IS-41, this local agent which the HLR points to can be replaced as the MU moves from one VLR to another. However, during its tenure period serving as a local agent for the MU, it executes the TLA scheme on behalf of the HLR.

Specifically, when the MU moves from the local agent to a nearby VLR, only a forwarding pointer is setup between the local agent and the new VLR, and no location update operation is performed to the HLR. We also refer to this cost as the VLR–VLR “binding” cost. After that, the MU can move back and forth between the two VLRs, including the local agent, without informing the local agent or the HLR. When a call arrives, the local agent will first attempt to locate the MU in its area. If the MU cannot be found, then it will follow the forwarding pointer to the next VLR to search for the called MU, thus increasing the search cost. Note that, however, this search cost is still less than that of the TLA scheme because on the second search attempt, the system only has to pay the VLR–VLR communication cost in the

Table 2
Databases stored in the HLR, local agent and MU under ATLA.

Movement	HLR pointer	Local agent pointer	MU table	Update cost
(a) $V_2 \rightarrow V_1$	V_2	V_1	(V_1^*, V_2)	VLR–VLR
(b) $V_1 \rightarrow V_2$	V_2	V_1	(V_2^*, V_1)	No
(c) $V_2 \rightarrow V_3$	V_2	V_3	(V_3^*, V_2)	VLR–VLR
(d) $V_3 \rightarrow V_4$	V_4	V_3	(V_4^*, V_3)	VLR–HLR

ATLA scheme by following the VLR–VLR link, as opposed to the VLR–HLR communication cost in the TLA scheme.

Note that in the ATLA scheme, the MU must also record two recently visited VLRs in the table as in TLA, including the local agent. We will call the non-agent VLR in the table as the 2nd VLR. When the MU moves from the local agent VLR to a new VLR, it will inform the local agent to update its forwarding pointer to point to the new VLR which it just entered, thus making the new VLR as the 2nd VLR, without notifying the HLR. However, if the move is from the 2nd VLR to a new VLR, then an update operation to the HLR must be performed, thus replacing the local agent with the new VLR just entered while keeping the 2nd VLR unchanged. Table 2 illustrates ATLA by showing the contents of the databases stored in the HLR and MU based on the ATLA scheme. After making movement (a), the HLR still records V_2 as the agent, the agent (i.e., V_2 at this point) sets up a pointer to V_1 , and the MU’s location table stores (V_1^*, V_2) , with V_1 being marked with “*” because it is the current VLR. In this case, the MU only sends a VLR–VLR binding message to the local agent without informing the HLR about the location change. After movement (b), the MU moves back to the previous VLR, V_2 , so it does not need to do anything except updating its own database to (V_2^*, V_1) . When the MU makes movement (c), it discovers that it moves into a new VLR V_3 , so it sends a VLR–VLR binding message to the local agent V_2 . Finally, when the MU moves into V_4 from the 2nd VLR V_3 , it performs an update operation to the HLR, thus replacing V_2 with V_4 as the new local agent.

Another note that is worth mentioning is the concept of state inconsistency. In table 2, we see that after movement (a) or (c), the system is in an inconsistent state, i.e., the pointer from the HLR points to a VLR that is not the current VLR, while after movement (b) or (d), the system is in a consistent state. This can also be observed from table 2 by combining columns 2 and 3 and see if it is the same as column 4 in sequence.

2.4. CATLA

The CATLA strategy is the same as the ATLA strategy except that an update operation to the HLR will occur whenever the MU moves across a network switch boundary. The idea is to reduce the average VLR–VLR communication cost between the local agent and the 2nd VLR. In this scheme, if the MU enters a new VLR covered by another network switch (e.g., a STP in figure 1), the MU will update the HLR directly because otherwise the VLR–VLR communication cost crossing the STP switch boundary is expensive and such

Table 3
Databases stored in the HLR, local agent and MU under CATLA.

Movement	HLR pointer	Local agent pointer	MU table	Update cost
(a) $V_2 \rightarrow V_1$	V_2	V_1	(V_1^*, V_2)	VLR-VLR
(b) $V_1 \rightarrow V_2$	V_2	V_1	(V_2^*, V_1)	No
(c) $V_2 \rightarrow V_3$	V_3	nil	(V_3^*, nil)	VLR-HLR
(d) $V_3 \rightarrow V_4$	V_3	V_4	(V_4^*, V_3)	VLR-VLR

cost would strike twice: once when the pointer is set up between the local agent and the 2nd VLR (which are under two different STPs), and once when a call arrives subsequently and a location miss results. Because of this force-update rule, we expect that CATLA may gradually perform better than ATLA as the CMR value of the MU increases. Table 3 illustrates CATLA by showing the contents of the databases maintained by the HLR, local agent, and MU based on CATLA for the same movement scenario discussed earlier. Movements (a) and (b) are the same as in the ATLA scheme, since no switch boundary is crossed. Movement (c), however, incurs a location update operation to the HLR since V_2 and V_3 are under two separate STPs (see figure 1), after which V_3 becomes the local agent and the 2nd VLR is nil. Movement (d) is also different from that in the ATLA scheme because V_3 now is the local agent, not the 2nd VLR as in the ATLA scheme. Therefore, when the MU enters V_4 from V_3 , only a pointer is setup from V_3 to V_4 since V_4 is under the same STP as V_3 .

3. Modeling the location strategies

In this section, we develop three separate Markov models to describe the performance characteristics of the PCS network operating under TLA, ATLA and CATLA. We will later apply the result obtained from these models to compare TLA, ATLA, CATLA and IS-41 under identical conditions based on a hexagonal network structure.

3.1. Modeling TLA

The state of a MU as it crosses database boundaries while being called can be described by a three-component state description vector (a, b, c) . Component a is a binary quantity indicating whether or not the mobile unit is in the state of being called, with 0 standing for idle and 1 standing for busy. Component b is also a binary quantity indicating if the MU has just moved to a new VLR area, with 1 meaning yes and 0 meaning no. The third component, c , indicates if the location table maintained by the mobile unit is inconsistent with that maintained by the HLR, with 0 meaning consistent and 1 meaning inconsistent.

Figure 2 shows the Markov model for describing the PCS operating under TLA. Table 4 shows the notation used for the TLA model. Initially, the MU is in state $(0, 0, 0)$, meaning that it is not being called, the MU has not yet made any move across any VLR boundary, and the location table stored in the HLR is consistent with that stored in the

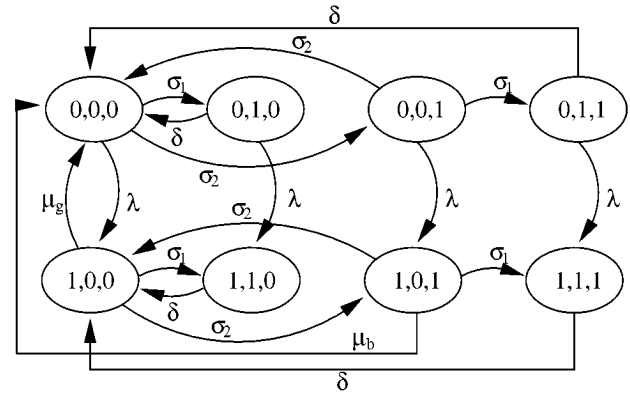


Figure 2. Markov model for PCS network under two-location algorithm.

Table 4
Notation used for the TLA model.

λ	arrival rate of calls to the MU
σ	mobility rate of the MU, i.e., the rate at which the MU crosses VLR boundaries
θ	probability of the MU moving back to the previous VLR
CMR	call-to-mobility ratio of the MU, i.e., $\text{CMR} = \lambda/\sigma$
σ_1	mobility rate of the MU moving to a new VLR, i.e., $\sigma_1 = (1 - \theta)\sigma$
σ_2	mobility rate of the MU moving to the previous VLR, i.e., $\sigma_2 = \theta\sigma$
δ	location update rate in updating the location table stored in the HLR
μ_g	location search rate in locating the MU when the location table in the MU is consistent with that in the HLR
μ_b	location search rate in locating the MU when the location table in the MU is inconsistent with that in the HLR

MU. Below, we explain briefly how we construct the Markov model.

First, if the MU is in state $(0, i, j)$, $0 \leq i, j \leq 1$, and a call arrives, then the new state is $(1, i, j)$, i.e., the MU is now in the state of being called. This behavior is modeled by the (downward) transition from state $(0, i, j)$ to state $(1, i, j)$, $0 \leq i, j \leq 1$, with a transition rate of λ .

Second, if the MU is in state $(1, i, j)$ and another call arrives, then the MU will remain at the same state, since the MU remains in the state of being called. This behavior is described by a hidden transition from state $(1, i, j)$ back to itself with a transition rate of λ . This type of transition is not shown in figure 2 since it does not need to be considered when solving a Markov chain [8]. Note that this implies that in state $(1, i, j)$ the number of requests accumulated to locate the MU may be greater than one.

Third, if the MU is in state $(1, 0, 0)$, it means that the location table stored in the HLR is consistent with that stored in the mobile unit and the mobile unit is in the state of being called. Therefore, the PCS network can service all pending calls simultaneously with a service rate of μ_g . After the service, the new state is $(0, 0, 0)$.

Fourth, if the MU is in state $(1, 0, 1)$, it means that the location table stored in the HLR is inconsistent with that stored in the mobile unit but there are pending calls waiting to be serviced. Therefore, the PCS network has to spend

twice as much time to locate the mobile unit. This behavior is modeled by using a different service rate of μ_b from state $(1, 0, 1)$ to state $(0, 0, 0)$. After the service, the new state is $(0, 0, 0)$ because the location table stored in the HLR is updated after the call delivery service and is therefore consistent with that stored in the mobile unit again.

Last, regardless of whether or not the MU is in the state of being called, the MU can move across a VLR boundary. There are two cases:

1. If the mobile unit moves to a new VLR then an update operation has to be performed to the table stored in the HLR. This behavior is modeled by a transition from state $(i, 0, j)$ to $(i, 1, j)$, $0 \leq i, j \leq 1$, with a transition rate σ_1 , after which the system transmits from state $(i, 1, j)$ to state $(i, 0, 0)$ with a transition rate of δ .
2. If the MU moves back to the previously visited VLR, then there is no update operation required to update the HLR associated with this event. This is modeled by a transition from state $(i, 0, 0)$ to state $(i, 0, 1)$, $0 \leq i \leq 1$, with rate σ_2 . This results in an inconsistent state, i.e., after the transition the location table stored in the HLR is consistent with that stored in the MU. Another possible transition is the reverse of the above, that is, from state $(i, 0, 1)$ to state $(i, 0, 0)$, $0 \leq i \leq 1$, also with a transition rate σ_2 . This, however, results in a consistent state. Note that the update time of the location table stored in the mobile unit with respect to the PCS network is zero (although the mobile unit's location table is updated itself). Therefore, there is no need to model the time needed to perform the update operation in this case.

The Markov chain shown in figure 2 is ergodic [8], which means that all states have a non-zero probability. The probability that the system is found in a particular state in equilibrium depends on the relative magnitude of the outgoing and incoming transitions rates. Let $P_{(i,j,k)}$ be the probability that the system stays in state (i, j, k) in equilibrium. Let TLA_{update} be the average cost of the PCS network in servicing a location update operation and let TLA_{call} be the average cost in locating the MU. Furthermore, let TLA_{cost} be the average cost of the PCS network in servicing the above two types of operations between two consecutive calls. Then,

$$TLA_{update} = \sum_{i=0}^1 (P_{(0,0,i)} + P_{(1,0,i)}) (1 - \theta) \frac{1}{\delta} + \sum_{i=0}^1 (P_{(0,1,i)} + P_{(1,1,i)}) \frac{1}{\delta}, \quad (1)$$

$$TLA_{call} = \sum_{i=0}^1 \sum_{j=0}^1 P_{(i,j,0)} \frac{1}{\mu_g} + \sum_{i=0}^1 \sum_{j=0}^1 P_{(i,j,1)} \frac{1}{\mu_b}, \quad (2)$$

$$TLA_{cost} = TLA_{update} \frac{\sigma}{\lambda} + TLA_{call}. \quad (3)$$

Equation (3) is obtained above because between two consecutive calls, the number of mobility moves across VLR

Table 5
Additional parameters used in the ATLA model.

δ_a	location update rate in updating the location table stored in the agent
μ_a	location search rate in locating the MU when the local agent is not the current VLR

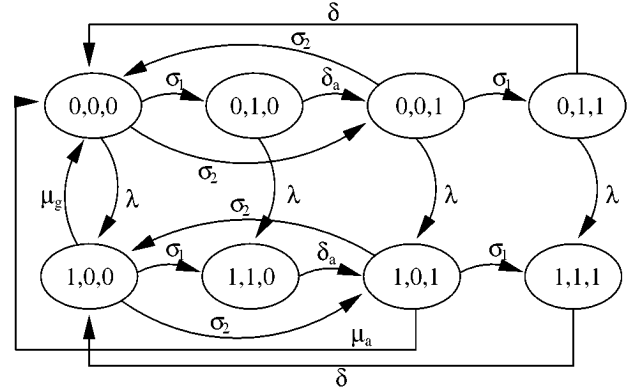


Figure 3. Markov model for PCS network under ATLA.

boundaries by the MU is equal to σ/λ on average. Note that the number of moves corresponds to the number of update operations, although some of which may not cause any update cost to the PCS network depending on whether or not the location table in the HLR needs to be updated.

3.2. Modeling ATLA

To model ATLA, we use the same set of parameters in table 4 without the last parameter μ_b , and introduce two more parameters listed in table 5. We describe ATLA by also using a three-component vector (a, b, c) . The meanings of components a and b are the same as before. The third component, c , indicates if the MU currently resides under the agent area, i.e., if the local agent is the current VLR, with 0 meaning that it is, and 1 meaning that it is not. Conceptually, the third component indicates if the HLR points to the current VLR, with 0 meaning yes and 1 meaning no. If yes, the system is in a consistent state in which the state information maintained by the HLR is consistent and the search operation can be done efficiently; otherwise, the system is in an inconsistent state and must follow a forwarding pointer to locate the MU.

Figure 3 shows the Markov model for describing the PCS operating under ATLA. Initially, the MU is in state $(0, 0, 0)$, meaning that it is not being called, has not moved across a VLR boundary, and the MU is currently located under the local agent (note that the HLR always points to the local agent in the ATLA scheme). Below, we explain briefly the differences between TLA and ATLA as we construct the Markov model for ATLA. The ATLA model is almost the same as the TLA model, except that we replace the transitions from state $(i, 1, 0)$ to state $(i, 0, 0)$ by those from state $(i, 1, 0)$ to state $(i, 0, 1)$, $0 \leq i \leq 1$, and replace μ_b by μ_a .

If the MU is in state $(1, 0, 1)$, it means that there are pending calls waiting to be serviced in this state but the MU is not in the cover area of the local agent, i.e., the local agent is not the current VLR. In this case, the system needs to follow the forwarding pointer from the local agent to locate the current VLR. This behavior is modeled by using a service rate of μ_a from state $(1, 0, 1)$ to state $(0, 0, 0)$. After the service, the new state is $(0, 0, 0)$ because the HLR database is updated after the call delivery service and is, therefore, consistent with that stored in the MU again. This rate μ_a is different from μ_b in figure 2 for the TLA model since only a VLR–VLR communication cost incurs in the second search attempt by following the forwarding pointer. Recall that for the TLA model, the second search attempt also involves a VLR–HLR communication cost.

When the MU moves across a VLR boundary, there are three cases to be considered:

1. If the MU moves from the local agent to a new 2nd VLR, then a binding operation has to be performed to bind the new 2nd VLR to the agent by means of a forwarding pointer. This behavior is modeled first by a transition from state $(i, 0, 0)$ to $(i, 1, 0)$, $0 \leq i \leq 1$, with a transition rate σ_1 , after which the system goes from state $(i, 1, 0)$ to state $(i, 0, 1)$ with a transition rate of δ_a . The first transition models the move event while the second models the binding (setting up pointer) event.
2. If the MU moves from a VLR, V_i , to a new VLR, V_j , where V_i and V_j are not the agent, then an update operation has to be performed to the table stored in the HLR. After the update operation, the location table of the HLR records the new VLR as the agent. This behavior is modeled by a transition from state $(i, 0, 1)$ to $(i, 1, 1)$, $0 \leq i \leq 1$, with a transition rate σ_1 , after which the system transits from state $(i, 1, 1)$ to state $(i, 0, 0)$ with a transition rate of δ . Note that the originating state in this case is $(i, 0, 1)$, $0 \leq i \leq 1$, with the 2nd state component being 0 (meaning it has not moved to a new VLR) and the 3rd state component being 1 (meaning the agent is not the current VLR).
3. If the MU moves back to the previously visited VLR, then there is no update operation required to perform the HLR associated with this update event. This is modeled by a transition from state $(i, 0, 0)$ to state $(i, 0, 1)$ with a transition rate σ_2 , after which the location table stored in the HLR is inconsistent with that stored in the MU, or from state $(i, 0, 1)$ to state $(i, 0, 0)$, $0 \leq i \leq 1$, also with a transition rate σ_2 , after which the location table stored in the HLR becomes consistent with that stored in the MU. Note that again the update time of the location table of the MU with respect to the PCS network is zero.

Let $ATLA_{\text{update}}$ be the average cost of the PCS network in servicing an update operation and let $ATLA_{\text{call}}$ be the average cost in locating the MU. Furthermore, let $ATLA_{\text{cost}}$ be the average cost of the PCS network in servicing the

above two types of operations between two consecutive calls. Then,

$$ATLA_{\text{update}} = \sum_{i=0}^1 P_{(i,0,0)}(1-\theta)\frac{1}{\delta_a} + \sum_{i=0}^1 P_{(i,1,0)}\frac{1}{\delta_a} + \sum_{i=0}^1 P_{(i,1,1)}\frac{1}{\delta} + \sum_{i=0}^1 P_{(i,0,1)}(1-\theta)\frac{1}{\delta}, \quad (4)$$

$$ATLA_{\text{call}} = \sum_{i=0}^1 P_{(i,0,0)}\frac{1}{\mu_g} + \sum_{i=0}^1 P_{(i,1,0)}\frac{1}{\mu_a} + \sum_{i=0}^1 \sum_{j=0}^1 P_{(i,j,1)}\frac{1}{\mu_a}, \quad (5)$$

$$ATLA_{\text{cost}} = ATLA_{\text{update}}\frac{\sigma}{\lambda} + ATLA_{\text{call}}. \quad (6)$$

3.3. Modeling CATLA

To model CATLA, we again use the same set of parameters in table 4 except μ_b , along with three more parameters, δ_c , μ_c and γ , listed in table 6. For notational convenience, let $\sigma_n = (1-\gamma)\sigma_1$ and $\sigma_c = \gamma\sigma_1$.

Figure 4 shows a Markov model for describing the PCS operating under CATLA. The same three-component state

Table 6
Additional parameters used in the CATLA model.

δ_c	location update rate in updating the location table of the local agent
μ_c	location search rate in locating the MU when the local agent is not the current VLR
γ	probability of the MU moving within the same network switch
σ_n	mobility rate of the MU moving to a new VLR not covered by the same network switch, i.e., $\sigma_n = (1-\gamma)(1-\theta)\sigma = (1-\gamma)\sigma_1$
σ_c	mobility rate of the MU moving to a new VLR covered by the same network switch, i.e., $\sigma_c = \gamma(1-\theta)\sigma = \gamma\sigma_1$

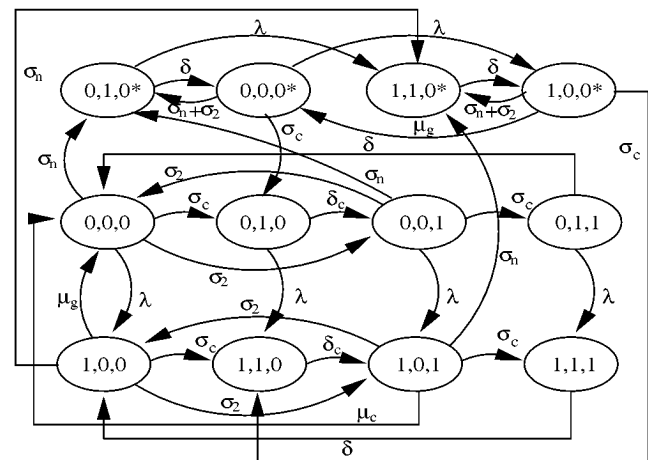


Figure 4. Markov model for PCS network under CATLA.

description as defined in the ATLA model is used here. Recall that under CATLA, whenever the MU moves across a network switch boundary, the HLR database is updated, in which case the new agent does not have a forwarding pointer in its local database, i.e., the 2nd VLR is nil. We differentiate this particular state by using the symbol “*”. For example, in state $(0, 0, 0)^*$ the forwarding pointer stored in the agent is nil, while in state $(0, 0, 0)$ the forwarding pointer does exist. Note that in both cases, the system is in a consistent state, i.e., the local agent is the current VLR. The Markov model for CATLA is self-explanatory. Below we discuss the behavior of the MU as it moves to a new VLR. There are three possible cases:

1. If the MU moves to a new VLR covered by a different network switch (the mobility rate of which is σ_n), then an update operation is performed to the database stored in the HLR. This behavior is modeled first by a transition from state $(i, 0, j)$ to $(i, 1, 0)^*$, $0 \leq i, j \leq 1$, with a transition rate σ_n , after which the system goes from $(i, 1, 0)^*$ to $(i, 0, 0)^*$ with a transition rate of δ to update the HLR’s database to point to the new VLR. After the update, the new VLR becomes the local agent of the MU, but in its database the forwarding pointer to the 2nd VLR is nil because in this case the 2nd VLR does not exist.
2. After the MU just crosses a network switch, it is possible that the MU may subsequently cross a network switch boundary again. This event comes in two forms: (a) the MU subsequently moves back to the previously visited VLR, the mobility rate of which is σ_2 ; (b) the MU simply enters another VLR not covered by the same switch, the mobility rate of which is σ_n . In either case, another update operation will be triggered to update the HLR database. This behavior is modeled first by a transition from state $(i, 0, 0)^*$ to state $(i, 1, 0)^*$, $0 \leq i \leq 1$, (the top row of figure 4) with a transition rate $\sigma_2 + \sigma_n$, after which the system goes from state $(i, 1, 0)^*$ to state $(i, 0, 0)^*$ with a transition rate of δ . After the MU just crosses a network switch, another likely event is that it subsequently goes to another VLR also within the same network switch. This is modeled first by a transition from state $(i, 0, 0)^*$ to state $(i, 1, 0)$, $0 \leq i \leq 1$, with a transition rate σ_c , after which the system goes from state $(i, 1, 0)$ to state $(i, 0, 1)$ with a transition rate of δ_c to update the pointer stored in the local agent. In this event, the local agent of the MU remains unchanged but the target state becomes inconsistent.
3. When the local agent and the 2nd VLR are not nil, if the MU moves to a VLR within the same network switch that covers the current VLR (which can be the local agent or the 2nd VLR), the system behaves the same as in the ATLA model. This behavior is described by the state transitions in the last two rows of figure 4. Note that here the third component in the state description tells us exactly whether or not the local agent is the current VLR. The target VLR which the MU moves into in

this case can be the previously visited VLR (the mobility rate of which is σ_2) or a new VLR (the mobility rate of which is σ_c). In both cases, the MU stays within the same network switch.

Again, let $CATLA_{update}$ be the average cost of the PCS network in servicing an update operation and let $CATLA_{call}$ be the average cost in locating the MU. Furthermore, let $CATLA_{cost}$ be the average cost of the PCS network in servicing the above two types of operations between two consecutive calls. Then,

$$\begin{aligned}
 CATLA_{update} &= \frac{1}{\delta_c} \left(\sum_{i=0}^1 (1-\theta)\gamma(P_{(i,0,0)} + P_{(i,0,0)^*}) + P_{(i,1,0)} \right) \\
 &+ \frac{1}{\delta} \sum_{i=0}^1 (1-\theta)(1-\gamma)P_{(i,0,0)} \\
 &+ \frac{1}{\delta} \left(\sum_{i=0}^1 (1-\theta)P_{(i,0,1)} + P_{(i,1,0)^*} + P_{(i,1,1)} \right) \\
 &+ [\theta + (1-\theta)(1-\gamma)]P_{(i,0,0)^*}, \quad (7)
 \end{aligned}$$

$$\begin{aligned}
 CATLA_{call} &= \sum_{i=0}^1 (P_{(i,0,0)} + P_{(i,1,0)^*} + P_{(i,0,0)^*}) \frac{1}{\mu_g} \\
 &+ \sum_{i=0}^1 (P_{(i,1,0)} + P_{(i,0,1)} + P_{(i,1,1)}) \frac{1}{\mu_c}, \quad (8)
 \end{aligned}$$

$$CATLA_{cost} = CATLA_{update} \frac{\sigma}{\lambda} + CATLA_{call}. \quad (9)$$

4. Application

In this section, we apply the results obtained in the last section to compare TLA, ATLA, CATLA and IS-41 under identical conditions. The network cost parameters listed in table 7 will apply to all location management schemes. Specific values of these network communication cost parameters can be obtained by considering specific network coverage models. We will show how to do so with a hexagonal coverage model. In this paper, we fix the (relative) value of U at 1.0 as used in [9], but use different values of T , τ and τ' to study their effects. We note that the relation among these last three parameters is

$$\tau = \gamma\tau' + (1-\gamma)T. \quad (10)$$

This is so because the STP switches are connected by the PSTN. Therefore, for any two VLRs not under the same STP, the VLR–VLR communication cost is comparable to the VLR–HLR communication cost. Here, γ is as defined in table 6.

Table 7
Communication cost parameters for all models.

U	the average cost for locating the MU under the basic scheme; it is also the locating cost under the TLA scheme when the HLR's location table is consistent with that stored in the mobile unit, as well as the locating cost under the ATLA and CATLA scheme to find the local agent.
T	the average VLR–HLR communication cost between HLR and VLR; it is also the update cost under the basic scheme and the cost to update the location table in the HLR under the TLA, ATLA and CATLA schemes. In figure 5, we illustrate the relationship between U and T as follows: $U = 2T + P_c,$ where P_c is the average paging cost needed to page a mobile user within a VLR.
τ	the average VLR–VLR communication cost for two VLRs in the network; it is also the cost of a binding operation under the ATLA scheme.
τ'	the average VLR–VLR communication cost for two VLRs under the same STP switch; it is also the cost of a binding operation under the CATLA scheme.

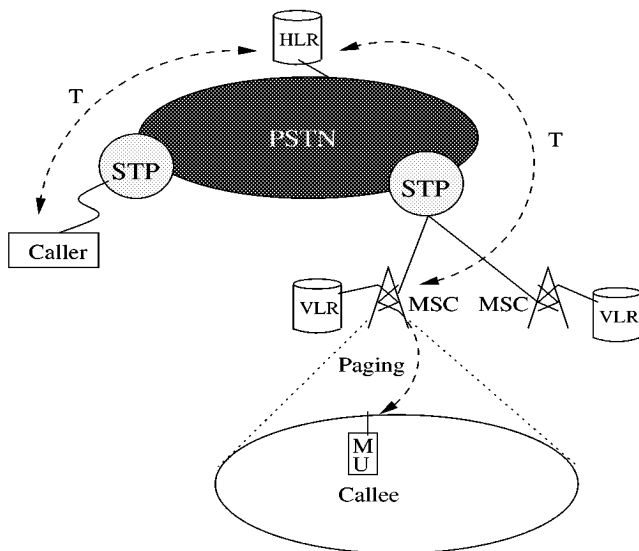


Figure 5. A search operation in the mobile network.

4.1. Parameterization of the TLA Markov models

There are six parameters in the TLA Markov model (see figure 2), i.e., σ , λ , θ , δ , μ_g , and μ_b (see table 4 for their meanings). Of these six parameters, σ , λ and θ are MU dependent parameters and will be studied in the paper by changing their values; on the other hand, δ , μ_g , and μ_b are network structure dependent and can be parameterized as

$$\delta = \frac{1}{T}, \quad \mu_g = \frac{1}{U}$$

and

$$\mu_b = \frac{1}{U + T + P_c} = \frac{1}{2U - T},$$

where U , T and P_c are as defined in table 7.

4.2. Parameterization of the ATLA Markov model

In ATLA, we again have three MU dependent parameters, i.e., σ , λ , and θ , which will be studied by changing their values. There are four network structure dependent parameters used in the ATLA Markov model (see figure 3), i.e., δ , μ_g , δ_a and μ_a (see tables 4 and 5 for their meanings). The parameterizations of δ and μ_g are the same as in TLA; δ_a and μ_a can be parameterized as follows:

$$\delta_a = \frac{1}{\tau}$$

and

$$\mu_a = \frac{1}{U + \tau + P_c} = \frac{1}{2U - 2T + \tau}.$$

4.3. Parameterization of the CATLA Markov model

In the CATLA Markov model (see figure 4), we have five network-structure-dependent parameters, i.e., δ , μ_g , δ_c , μ_c and γ (see tables 4 and 6 for their meanings) and the same three MU dependent parameters (σ , λ , and θ). Of these five network structure dependent parameters, δ and μ_g can be parameterized the same as before; γ can be derived given a specific network coverage model for which we will show how it can be done with a hexagonal coverage model; δ_c and μ_c can be parameterized as follows:

$$\delta_c = \frac{1}{\tau'}$$

and

$$\mu_c = \frac{1}{U + \tau' + P_c} = \frac{1}{2U - 2T + \tau'}.$$

4.4. Comparing TLA, ATLA, CATLA and IS-41

We use IS-41 as the baseline model against which our proposed agent-based forwarding strategies will be compared. The average network cost of the PCS network for location management under IS-41 is given by

$$\text{IS-41}_{\text{cost}} = \text{IS-41}_{\text{update}} \frac{\sigma}{\lambda} + \text{IS-41}_{\text{call}}, \quad (11)$$

where $\text{IS-41}_{\text{update}} = T$ and $\text{IS-41}_{\text{call}} = U$, as defined in table 7.

As in [7], we first study a case in which the VLR–VLR communication cost is one half of the VLR–HLR communication cost, i.e., $\tau/T = 0.5$. To derive γ , we consider a hexagonal network structure model wherein the number of VLRs under a n -layer STP is given by $3n^2 - 3n + 1$. Figures 6(a) and (b) show the number of VLRs covered by a 3-layer STP and a 2-layer STP, respectively, based on this hexagonal network coverage model.

Here we consider the case when $n = 2$, i.e., a STP covers 7 VLRs; later, we will study other cases to analyze the effect of n . For the case when $n = 2$, it can be shown that $\gamma = 0.57$ (see [2]). Therefore, when τ/T is 0.5, τ'/T is 0.125 based on equation (10). Figure 7 shows the average cost of the PCS

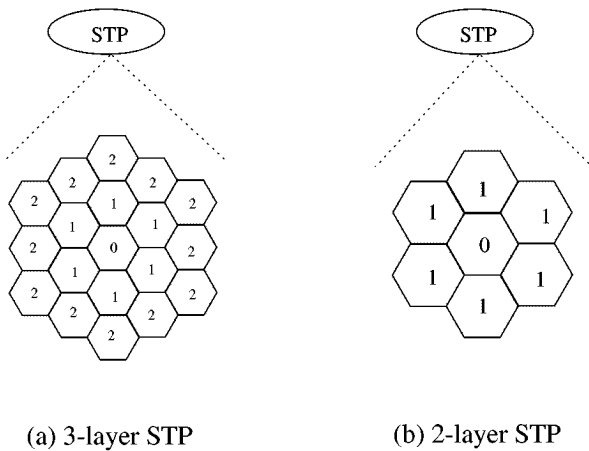


Figure 6. n -layer STPs under the hexagonal network coverage model.

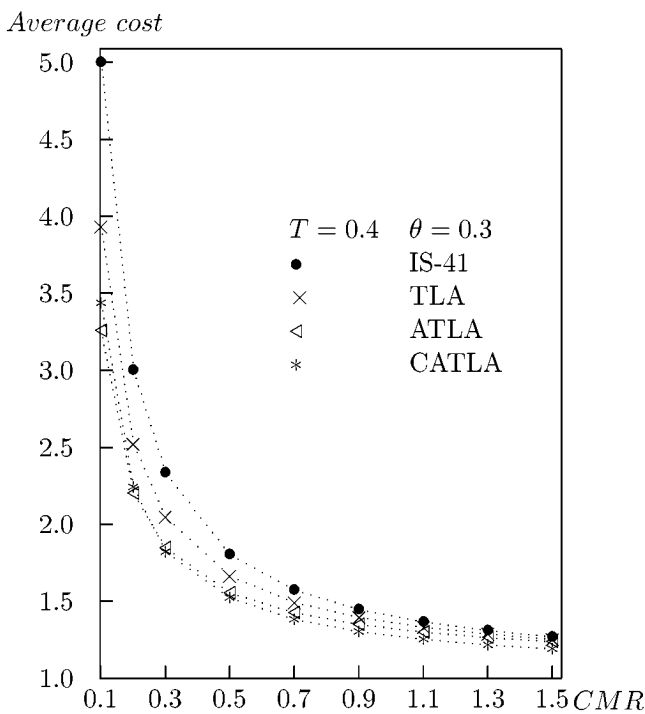


Figure 7. Comparison of IS-41, TLA, ATLA and CATLA with $T = 0.4$ and $\theta = 0.3$.

network due to location management under TLA, ATLA, CATLA and IS-41, as a function of CMR ($CMR = \lambda/\sigma$), for the case when $T = 0.4$ and $\theta = 0.3$ (with U fixed at 1). The data on the diagram were obtained by first solving a Markov model using the SHARPE software package [12] to obtain $P_{(i,j,k)}$ for each state (i, j, k) and subsequently computing the average cost based on equations developed in the paper (e.g., equation (3) for the TLA scheme). Figure 7 illustrates that when the CMR value is small, TLA, ATLA and CATLA all can outperform IS-41, even with θ as small as 0.3. For example, when CMR is 0.1, TLA, ATLA and CATLA schemes reduce the cost of IS-41 by 22%, 35% and 31%, respectively. In addition, ATLA and CATLA perform significantly better than TLA. CATLA is worse than ATLA when CMR is small, but is better as CMR's value increases.

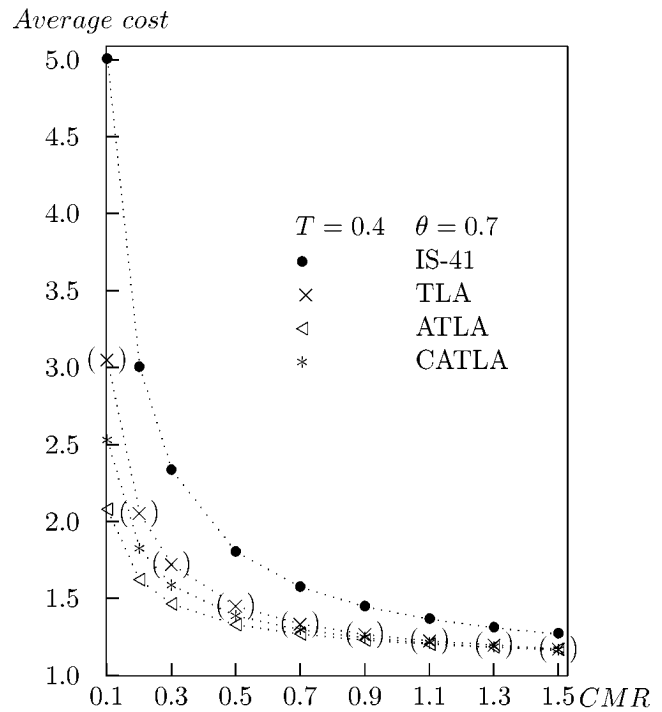


Figure 8. Comparison of IS-41, TLA, ATLA and CATLA with $T = 0.4$ and $\theta = 0.7$.

Here one should note that the performance improvement of ATLA and CATLA over TLA is significant because the cost metric used only accounts for the location cost between two consecutive calls, so the cumulative effect will be significant over the lifetime of the MU.

Figure 8 shows a similar condition as in figure 7 except that $\theta = 0.7$. In this particular case, TLA, ATLA and CATLA significantly outperform IS-41 and also their own respective counterparts in figure 7 with $\theta = 0.3$. This is because the probability of the MU moving back to the previously visited VLR is high ($\theta = 0.7$), and hence, the saving in the update cost is significant. This is especially the case when CMR is small partly because it costs less to update the user's location (since most moves are local) and partly because call arrivals are rare, so the high search cost can be amortized. Figure 8 also shows that in this case ATLA and CATLA can perform much better than TLA, because the VLR-VLR binding cost is less than the VLR-HLR update cost. Also, in this case since there is a high possibility of the MU moving back and forth, the advantage of CATLA over ATLA disappears even at high CMR values. This is conceivable since CATLA involves the VLR-HLR communication cost whenever a switch boundary is crossed and it is possible that the two most recently visited VLRs by the MU are covered by two separate network switches.

Both figures 7 and 8 are based on the condition that T is 0.4 of U , that is, the VLR-HLR communication cost is a large fraction of U (recall that $U = 2T + P_c$). In figure 9, we show an extreme case in which $T = 0.2$ while all other parameters remain the same as in figure 7. In this extreme case, since T is small compared to U , the advantage of TLA, ATLA and CATLA over IS-41 is not significant. In fact, the

advantage exists only at low CMR values. This is due to two reasons. First, at $\theta = 0.3$, the benefit of the local agent running the two-location algorithm on behalf of the HLR is not pronounced because the local agent will be replaced frequently, thus causing a high location update cost. Second, at $T = 0.2$, the HLR/VLR communication cost is relatively low, so there is little advantage for these algorithms to reduce the VLR–HLR communication cost. That is, whenever the MU moves into a new VLR, the location table of the HLR will be updated and the update cost becomes more affordable with a low T value. The benefit of this more frequent update operation is that it can keep a more precise view of the current VLR. Consequently, when a call arrives, it can serve the call with a lower cost. When CMR is larger than 0.5, the saving in update costs in ATLA and CATLA has a price. That is, ATLA and CATLA maintain a less precise state information than IS-41 and TLA, e.g., when the MU moves into a new VLR from the local agent, ATLA and CATLA do not update the HLR database and, thus, can enter an inconsistent state more often than TLA. As a result, when calls do arrive frequently, ATLA and CATLA must pay a higher search cost.

The classic tradeoff between the update and search costs is affected significantly by the MU's CMR. Figure 9 shows that at low CMR values, ATLA and CATLA perform better than TLA but the converse is true at high CMR values. In fact, figure 9 also shows that at relatively large CMR values, the basic IS-41 HLR/VLR scheme performs the best among all. The same physical interpretation regarding the tradeoff between the search and update costs applies. This, however, occurs only under extreme conditions when the probability

of going back and forth is low ($\theta = 0.3$) and the VLR–HLR communication cost ($T = 0.2U$) is a small fraction of U . For most PCS networks, we do expect that the VLR–HLR communication cost accounts for a large percentage of U (as in figure 8), and therefore, our proposed ATLA and CATLA schemes will outperform both TLA and IS-41. An important conclusion from our result is that at low CMR values, ATLA and CATLA always perform better than TLA.

4.5. Effect of network structure

Below we discuss the effect of network structures on the performance of agent-based forwarding strategies. Again suppose that the network structure is the hexagonal coverage model. Each STP is like a n -layer hierarchical structure covering $3n^2 - 3n + 1$ VLRs, where n can be either 2 or 3. It can be shown that the probability of the MU moving within the same STP (i.e., the γ parameter in table 6) is equal to 0.57 for $n = 2$ and 0.74 for $n = 3$ [2]. Furthermore, as n becomes larger, the average VLR–VLR communication cost (i.e., τ) becomes lower relative to the VLR–HLR communication cost (i.e., T) since most MU movements will be likely to be under the same STP, thus lowering the binding cost in the ATLA and CATLA schemes. The price to pay is to use more sophisticated STP switches to cover more VLRs under the same switch. The switch structure will affect TLA to a smaller extent because the second search attempt under TLA must always involve a VLR–HLR communication cost.

Below we consider a case when $T = 0.4$ and $\theta = 0.7$ as in figure 8. For $n = 2$, we again have $\tau/T = 0.5$ and $\tau'/T = 0.125$ as derived earlier. For $n = 3$, however, the value of τ is decreased due to the fact that a STP now covers

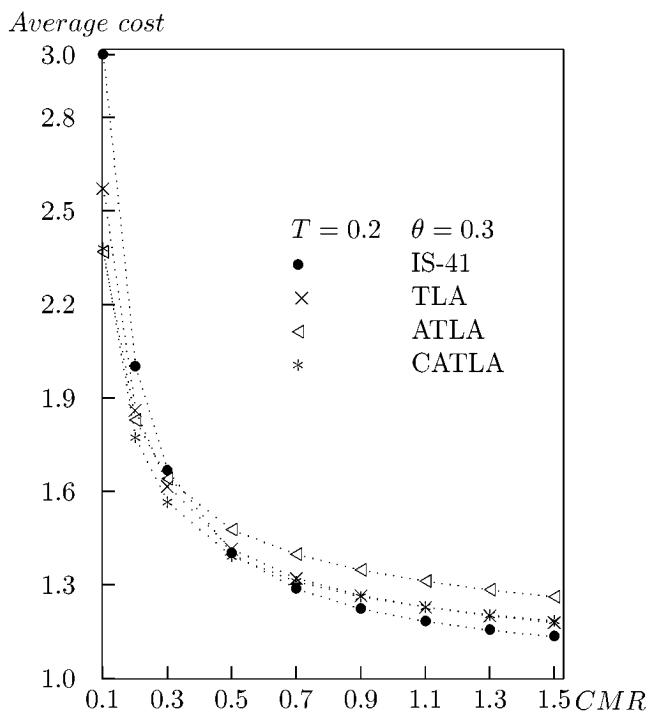


Figure 9. Comparison of IS-41, TLA, ATLA and CATLA with $T = 0.2$ and $\theta = 0.3$.

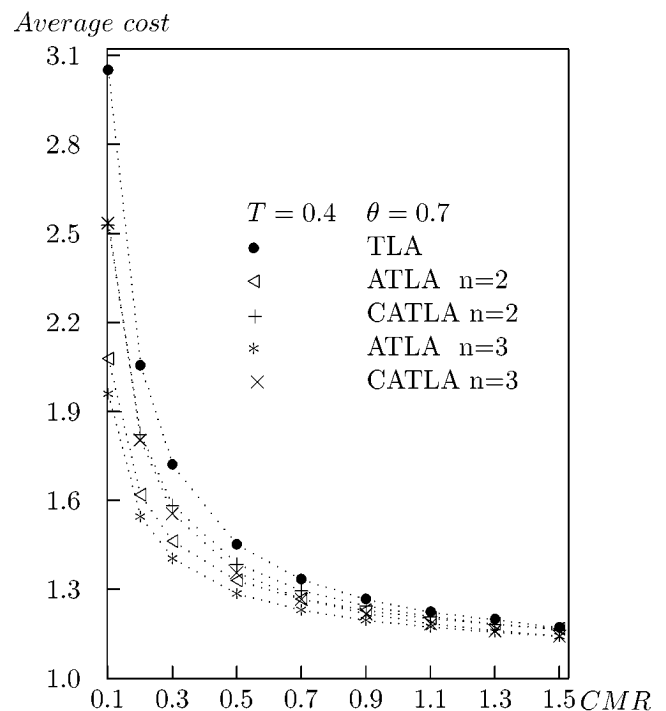


Figure 10. The effect of n -layer network structure on TLA, ATLA and CATLA.

more VLRs. To study the effect of n , let us fix $\tau'/T = 0.125$ since any two VLRs under the same STP will have the same VLR–VLR communication cost regardless of the number of VLRs under a STP. Therefore, based on equation (10), τ/T is 0.36 when $n = 3$.

Figure 10 shows the effect of n by comparing our proposed agent-based forwarding schemes against TLA for $n = 2$ and $n = 3$ for the case when $T = 0.4$ and $\theta = 0.7$. We see that as n increases, the performance gain of our proposed agent-based forwarding schemes over TLA becomes more and more pronounced as a result of a lower and lower VLR–VLR communication cost relative to the VLR–HLR communication cost. This means that our proposed schemes are especially attractive for network structures using sophisticated STPs that can cover many VLRs under one network switch.

5. Conclusion and future work

In this paper, we proposed the concept of agent-based forwarding to reduce the location management cost in mobile networks. In particular, we developed the ATLA and CATLA schemes which extend and improve over the TLA scheme. Markov models were used to analyze the performance characteristics of these algorithms. The exact condition under which one scheme is superior to the others and by how much can be assessed by using our Markov models. Our performance analysis data showed that our proposed schemes consistently outperform the basic IS-41 and TLA schemes for the same condition, especially when the CMR value of the MU is low. The advantage of our approach is especially pronounced for a hierarchically structured network in which the VLR–VLR communication cost is much lower than the VLR–HLR communication cost by means of sophisticated STPs which can cover many VLRs under one switch.

Some future research areas related to this paper include (a) investigating the possibility of combining the agent-based forwarding scheme with other approaches such as caching to further reduce the network location management cost for a wider range of parameter values, e.g., covering both high and low CMR values; (b) designing “service hand-off” management algorithms in mobile systems and applying similar modeling techniques to assess their performance characteristics.

References

- [1] G. Cho and L.F. Marshall, An efficient location and routing scheme for mobile computing environments, *IEEE Journal on Selected Areas in Communications* 13(5) (June 1995) 868–879.
- [2] I.R. Chen, T.M. Chen and C. Lee, Performance evaluation of forwarding strategies for location management in mobile networks, *The Computer Journal* 41(4) (August 1998) 243–253.
- [3] EIA/TIA, Cellular radio telecommunication inter system operations, Technical report IS-41 (Revision B) (July 1991).
- [4] J.S.M. Ho and I.F. Akyildiz, Dynamic hierarchical database architecture for location management in PCS networks, *IEEE/ACM Transactions on Networking* 5(5) (October 1997) 646–660.
- [5] B. Jabbari, G. Golombo, A. Nakajima and J. Kulkarni, Network issues for wireless communications, *IEEE Communication Magazine* (January 1995).
- [6] R. Jain, Y.B. Lin, C. Lo and S. Mohan, A caching strategy to reduce network impacts of PCS, *IEEE Journal on Selected Areas in Communications* 12(8) (October 1994) 1434–1444.
- [7] R. Jain, Y.B. Lin, C. Lo and S. Mohan, A forwarding strategy to reduce network impacts of PCS, in: *14th Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE INFOCOM '95)*, Vol. 2 (1995) pp. 481–489.
- [8] L. Kleinrock, *Queueing Systems, Vol. 1, Theory* (Wiley, 1975).
- [9] Y.B. Lin, Reducing location update cost in a PCS network, *IEEE/ACM Transactions on Networking* 5(1) (February 1997) 25–33.
- [10] P. Lin and W.-N. Tsai, A simulation study of multiple-location scheme for PCS mobility management, in: *3rd Workshop on Mobile Computing*, Taiwan (1997) pp. 37–45.
- [11] M. Mouly and M.B. Pautet, *The GSM System for Mobile Communications* (M. Mouly, Palaiseau, France, 1992).
- [12] R. Sahner, K.S. Trivedi and A. Puliafito, *Performance and Reliability Analysis of Computer Systems: An Example-based Approach Using the SHARPE Software Package* (Kluwer Academic, 1996).
- [13] S. Tabbane, An alternative strategy for location tracking, *IEEE Journal on Selected Areas in Communications* 13(5) (June 1995).

Ing-Ray Chen received the BS degree from the National Taiwan University, Taipei, Taiwan, in 1978, and the MS and PhD degrees in computer science from the University of Houston, University Park, Houston, TX, in 1985 and 1988, respectively. He is currently an Associate Professor in the Department of Computer Science, Virginia Polytechnic Institute and State University. His research interests are in reliability and performance analysis, mobile computing, multimedia, and distributed systems. Dr. Chen serves on the Editorial Board of *IEEE Transactions on Knowledge and Data Engineering* and *The Computer Journal*. He is a member of the IEEE/CS and ACM.

E-mail: irchen@cs.vt.edu

Tsong-Min Chen received his B.A. degree from the National Chung-Hsing University in applied mathematics, Taichung, Taiwan, and the M.S. degree in computer science from the National Tsing-Hau University at Hsinchu in 1989 and 1991, respectively. Currently, he is a Ph.D. student in the Department of Computer Science and Information Engineering at the National Cheng Kung University. His research interests include performance modeling and mobile computing.

E-mail: tmchen@dbserver.csie.ncku.edu.tw

Chiang Lee received the B.S. degree from the National Cheng Kung University, Taiwan, in 1980 and the M.E. and Ph.D. degrees in electrical engineering from the University of Florida, Gainesville, Florida, in 1986 and 1989, respectively. He joined IBM Mid-Hudson Laboratories, Kingston, NY in 1989 and participated in a project working on the design and performance analysis of a parallel and distributed database system. He joined the faculty of National Cheng Kung University in 1990 and is currently a Professor of Institute of Information Engineering of the University. His research interests are in the areas of mobile computing, web information systems, and integration of databases. He has many papers published in major database journals and conferences, and has been invited as an author of a chapter for several technical books. Dr. Lee has served as a Steering Committee member of the DASFAA International Conference from 1996 to 1998. He has also served on organizing and program committees for several major international conferences, including the IEEE International Conference on Data Engineering, the International Conference on Very Large Data Bases, and others.

E-mail: leec@csie.ncku.edu.tw, leec@dbserver.csie.ncku.edu.tw