

Consistency Mechanisms for a Distributed Lookup Service supporting Mobile Applications

Christoph Lindemann and Oliver P. Waldhorst

University of Dortmund

Department of Computer Science

August-Schmidt-Str. 12

44227 Dortmund, Germany

<http://www4.cs.uni-dortmund.de/~Lindemann/>

ABSTRACT

This paper presents a general-purpose distributed lookup service, denoted **Passive Distributed Indexing (PDI)**. PDI stores entries in form of (key, value) pairs in index caches located in each mobile device. Index caches are filled by epidemic dissemination of popular index entries. By exploiting node mobility, PDI can resolve most queries locally without sending messages outside the radio coverage of the inquiring node. Thus, PDI reduces network traffic for the resolution of keys to values. For keeping index caches coherent, configurable value timeouts implementing implicit invalidation and lazy invalidation caches implementing explicit invalidation are introduced. Inconsistency in index caches due to weak connectivity or node failure is handled by value timeouts. Lazy invalidation caches reduce the fraction of stale index entries due to modified data at the origin node. Similar to index caches, invalidation caches are filled by epidemic distributions of invalidation messages. Simulation results show that with the suitable integration of both invalidation mechanisms, more than 95% of results delivered by PDI index caches are up-to-date for the considered scenario.

Categories and Subject Descriptors

C.2.4 [Computer Communication Networks]: Distributed Systems – *Distributed Applications, Distributed Databases*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval – *Search process*; C.4 [Performance of Systems] – *Design Studies*;

General Terms

Algorithms, Performance, Design, Experimentation.

Keywords

Consistency maintenance in weakly connected environments, data caching, ad-hoc networked databases, epidemic algorithms for data dissemination.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

MobiDE'03, September 19, 2003, San Diego, California, USA.

Copyright 2003 ACM 1-58113-767-2/03/0009...\$5.00.

1. INTRODUCTION

Many distributed applications require global resolution of application-specific keys to application-specific values, a functionality provided by a lookup service. Perhaps the most prominent lookup service is the Domain Name System (DNS, [11]), which resolves each host name to the corresponding IP address. Similarly, instant messaging systems require a one-to-one mapping of a user ID to the current user's terminal or availability state. Further examples of lookup services include Internet search engines and distributed information retrieval systems, which provide for each query consisting of keywords a one-to-many mapping to matching documents. Recent research efforts in peer-to-peer technology [14], [15] aim at building Internet-scale distributed hash tables which provide a general-purpose approach for mapping keys to values.

In mobile and wireless environments, weak connectivity or even disconnected operation hampers the employment of a centralized lookup service. The smart collaboration of mobile devices in an ad hoc fashion constitutes an attractive alternative for implementing an effective distributed lookup service for such scenarios. Suitable methods for implementing such collaboration constitute epidemic algorithms. Such algorithms transmit information when nodes get in direct contact, similar to the transmission of an infectious disease between individuals. Mathematical models for the spread of epidemic diseases have been widely studied. There exist applications of epidemic algorithms in various fields of computer science, e.g., for the maintenance of replicated databases [2]. Papadopoulou and Schulzrinne introduced seven degrees of separation (7DS), a system for mobile Internet access based on web document dissemination between mobile users [12], [13]. To locate a web document, a 7DS node broadcasts a query message to all mobile nodes currently located inside its radio coverage. Recipients of the query send response messages that contain file descriptors of matching web documents stored in their local file caches. Subsequently, such documents can be downloaded with HTTP by the inquiring mobile node. Web documents may be distributed to other nodes that move into radio coverage, implementing an epidemic spread of information.

Using a related approach, Goel, Singh, Xu and Li proposed broadcasting segments of shared files using redundant tornado encoding [4]. Their approach enables nodes to restore a file, if a sufficient number of different segments have been received from one or more sources. Khelil, Becker, Tian, and Rothermel

presented an epidemic model for a simple information diffusion algorithm [8] inspired by the SPIN-1 protocol [5]. Both systems implement a push model for information dissemination. That is, shared data is advertised or even actively broadcasted without a node requesting it. Recently, Hanna, Levine, and Mamatha proposed a fault-tolerant distributed information retrieval system for peer-to-peer document sharing in mobile ad hoc environments [6]. Their approach distributes the index of a new document to a random set of nodes when the document is added to the system. The complete index of a document, i.e., all keywords matching it, constitutes the smallest unit of disseminated information.

In this paper, we introduce a general-purpose distributed lookup service, denoted Passive Distributed Indexing (PDI). Building upon [9], PDI stores index entries in form of (key, value) pairs in index caches located in each mobile device. Index caches are filled by epidemic dissemination of popular index entries. Thus, by exploiting node mobility PDI can resolve most queries locally without sending messages outside the radio coverage of the inquiring node. Thus, PDI effectively reduces network traffic for the resolution of keys to values for applications possessing a sufficiently high degree of temporal locality in their query streams as known for web content, Internet search engines and P2P file sharing systems [16], [18]. Beyond [9], we generalize the semantics of key-to-value matching beyond the simple matching of keywords to document-sources. Thus, PDI supports arbitrary application-specific keys to be matched to application-specific values. For keeping index caches coherent, we introduce two novel consistency mechanisms for PDI index caches: (1) configurable value timeouts implementing implicit invalidation and (2) lazy invalidation caches implementing explicit invalidation. Inconsistencies in index caches due to weak connectivity or node failure are handled by value timeouts. Lazy invalidation caches reduce the fraction of stale index entries due to modified data at the origin node. Similar to the epidemic distribution of index entries, invalidation caches are filled by epidemic distributions of invalidation messages.

As for PDI, a sufficiently high degree of locality in user queries is essential for all related approaches [4], [6], [8], [12], [13]. Opposed to 7DS [12], [13], PDI implements document search based on a distributed index rather than immediate document sharing. In fact, PDI disseminates individual index entries where as 7DS disseminates complete documents. Opposed to [4] and [8],

PDI implements a pull model and does not advertise data yielding a substantial reduction of traffic over wireless links. In contrast to [6], PDI disseminates just a single key and one or more matching values rather than the complete index of a document. Furthermore, opposed to all related work, PDI provides effective means for coping with stale index entries due to weak connectivity, node failure, and modified data. Moreover, opposed to all related work, PDI cannot only be employed for P2P file sharing, but also for numerous other mobile applications like mobile instant messaging or a mobile information portal.

In a comprehensive simulation study using ns-2, [3], we investigate the coherence of the key-to-value matching provided by PDI. The presented performance curves show that value timeouts reduce inconsistent results delivered by index cache due to intermitted connectivity or node failure by more than 75%. The employment of lazy invalidation caches reduces the fraction of stale index results by more than 50%. With the integration of both invalidation mechanisms, more than 95% of the results delivered by PDI index caches are up-to-date for the considered scenario.

This paper is organized as follows. Section 2 summarizes the basic concepts of PDI. Section 3 introduces the two novel consistency mechanisms for PDI, value timeouts and lazy invalidation caches. Performance results for PDI with these consistency mechanisms are presented in Section 4. Finally, concluding remarks are given.

2. Basic Concepts of PDI

We assume a system consisting of several mobile nodes, e.g. mobile users equipped with notebooks or PDAs and wireless network interfaces as illustrated in Figure 1. All mobile nodes collaborate via a shared application that uses a distributed lookup service. Radio coverage is small compared to the area covered by all nodes, so that most nodes cannot contact each other directly. Thus, communication may be performed using several intermediate hops as in mobile ad hoc networks (MANET, [10]). Subsequently, we assume IEEE 802.11x as underlying radio technology [7]. However, we would like to point out that PDI could be employed on any radio technology that enables broadcast transmissions inside a node's radio coverage.

PDI implements a general-purpose lookup service for mobile applications. In general, PDI stores index entries in the form of

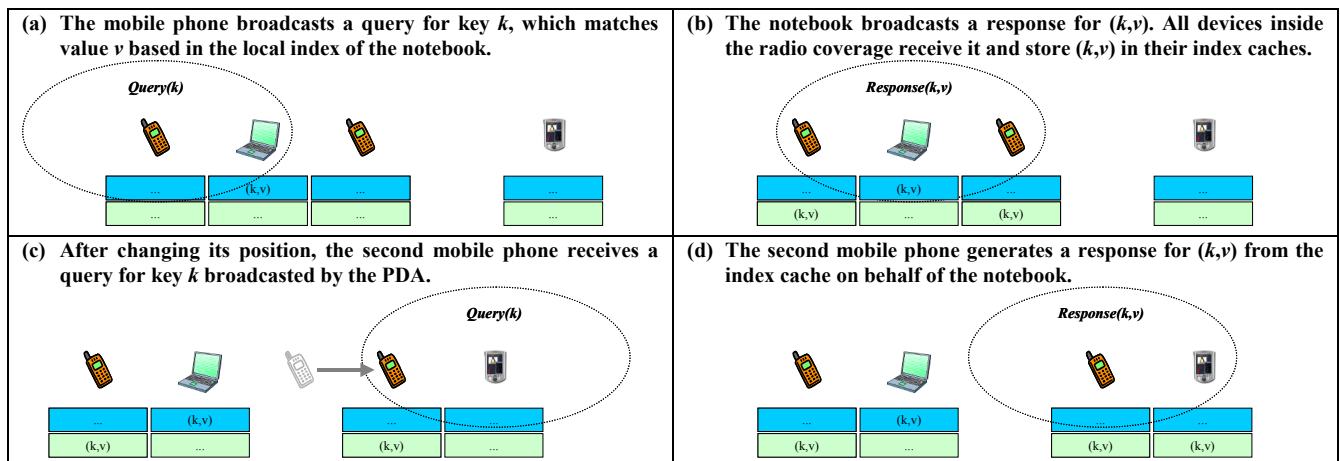


Figure 1. Illustration of epidemic information dissemination with PDI

pairs (k, v) . Keys k and values v are both defined by the mobile application. In case of file sharing, keys are given by keywords derived from the file name or associated meta data. Values are given by references to files in form of URIs. Opposed to distributed hash-table systems [14], [15], PDI does neither limit the number of keys matching a value nor the number of values matched by a key. However, some mechanisms introduced in Section 3 require that a value is unique in the system, i.e., it is only added to the system by a single node. This can be easily achieved by extending the application specific value v by a unique node identifier i_n for node n . For example, the node identifier i_n may be derived from the node's IP address or the MAC address of the radio interface. For ease of exposition, we will abbreviate the unique value given by (v, i_n) pairs just by v .

A node n may contribute index entries of the form (k, v) to the system by inserting them in a local index. In Figure 1, the local index is drawn as the first box below each mobile device. We refer to such an index entry as supplied. The node n is called the origin node of an index entry. For example, the notebook shown in Figure 1 is the origin node of the index entry (k, v) . A key k matches a value v , if (k, v) is currently supplied to the PDI system. Each node in the system may issue queries in order to resolve a key k to all matching values v_i (see Figure 1a). We refer to a node issuing a query as the inquiring node. We restrict ourselves to a query semantic given by conjunctions of keys, i.e., implementing Boolean AND. That is, for a query comprising of more than one key, a value matches the query, if it matches all keys. Note that the query semantic can be easily extended to additionally support disjunctions of keys, i.e., implementing Boolean OR, by including a bit vector indicating the matched keys for each value in the response message.

Query messages are sent to the IP limited broadcast address 255.255.255.255 and a well-defined port. Thus, all nodes located inside the radio coverage of the inquiring node receive a query message. These nodes may generate a response message. A response message contains the query and all matching values from either the local index or a second data structure called index cache. To enable epidemic data dissemination, PDI response messages are sent to the IP limited broadcast address 255.255.255.255 and a well-defined port, too. Thus, all mobile nodes within the radio coverage of the responding node will overhear the message (Figure 1b). Not only the inquiring node but also all other mobile nodes that receive a response message extract all index entries and store them in the index cache (see Figure 1b). In Figure 1, index caches are drawn as the second box below mobile devices. Index entries from the index cache are used to locally resolve queries, if the origin nodes of matching values reside outside the radio coverage of the inquiring node (see Figures 1c and 1d). Obviously, the index cache size is limited to a maximum number of entries adjusted to the capabilities of the mobile device. The replacement policy least-recently-used (LRU) is employed when newly received index entry will see a full cache. Note that information is disseminated to all other nodes that are in direct contact, similar to the spread of an infectious disease. Due to the movement of nodes and overhearing response messages of neighboring nodes, index entries are disseminated within the network without costly global communication. In fact, PDI builds and maintains an index distributed among mobile node of the MANET in a passive way.

Recall that responses derived from the index caches only contain index entries from remote hosts. Such responses may be out of date when the origin node either has left the system or has withdrawn the index entry from its local index. Thus, effective consistency mechanisms as introduced in Section 3 are the key for the effective design of PDI.

To extend information dissemination beyond the radio coverage of inquiring nodes, PDI includes a message forwarding mechanism. Queries may be relayed a certain number of hops specified by the inquiring node in a time-to-live field (TTL_{query}). Similarly, response messages may be forwarded TTL_{query} hops. Before forwarding a response message, a mobile node removes all values found in the index cache from the message. For a detailed description of this concept, denoted as selective forwarding, we refer to [9].

Note that beside inconsistent results, PDI might even return no results to a query for key k at all. This may occur, if the corresponding query neither reaches the origin node nor another node storing an index entry for value v in its index cache. We refer to such unresolved query as false misses. For the application running on top of PDI, there are two ways to resolve false misses: First, the application can tolerate false misses and simply ignore them. Second, the application has to resort to a centralized lookup service via the cellular infrastructure of a mobile network as fallback. Since query streams for P2P file sharing and Internet search engines possess a high degree of temporal locality [16], [18], PDI will produce only few false misses for such applications.

3. Means for Maintaining Index Consistency

3.1 Consistency Issues for Mobile Systems

As mentioned above, caching index entries may introduce inconsistency. For instance, when a node stops supplying an index entry (i.e., the pair (k, v) is removed from the local index of the node), copies of the index entry will remain in the index caches of other nodes. We refer to an index entry contained in a response message as fresh or up-to-date, if it is currently stored in the local index of the origin node. Otherwise, the index entry is denoted as stale. Stale index entries obviously yield stale search results and even disseminate in the system. As shown in Section 4.2, hits for stale index entries constitute a significant fraction of all results received in response to a query when no invalidation mechanism is used. In this section, we discuss invalidation mechanisms that can cope with both weak connectivity and modification of information. Since PDI has to deal with both sources of inconsistency, PDI implements an integrated approach that combines both invalidation mechanisms.

Timeouts constitute a common concept in several areas of distributed applications, as they can assure cache consistency without the need to contact the source of the cached information. Examples include the invalidation of cached DNS records in the domain name system or the invalidation of cached Web documents in WWW caches. To achieve both maximum consistency and a sufficient number cache hits, it is crucial for a timeout-based invalidation mechanism that the timeout durations are chosen appropriately. In DNS, the origin DNS server specifies the timeout duration for each entry. In web caching, the adaptive TTL algorithm calculates the timeout duration from the last modification time of a document at the origin server. Note that

both approaches rely on information directly received from the origin server. In contrast, due to the epidemic dissemination, most index cache entries are extracted from responses comprising of entries of other index caches, i.e., without direct contact to the origin node. Thus, PDI defines the concept of value timeouts to approximate the most recent information about the state of an index entry at the origin node.

Examples of explicit invalidation schemes include the invalidation of cached memory blocks in distributed shared memory (DSM) systems, or the invalidation of documents in web caches. To achieve consistency, the origin node of an item sends invalidation messages to exactly those nodes that are caching this item. In DSM systems, the origin node of a shared page sends invalidation messages to all nodes sharing this page. In web caching systems, the origin server of a web document sends invalidation messages to each web cache that holds a copy of the document. Note that both mechanisms require that the origin node knows where all copies of an item reside and that all sharers are reachable. In contrast, in a mobile environment consisting of nodes with limited resources, connectivity of nodes cannot be guaranteed nor directories for all cached copies of a shared item can be maintained. To address these constraints in mobile systems, PDI defines the concept of lazy invalidation caches implementing explicit invalidation of values.

3.2 Configurable Value Timeouts for Dealing with Weak Connectivity and Node Failures

The basic concepts of PDI as described in Section 2 do not take into account low connectivity and spontaneous departures of nodes; circumstances under which all information previously supplied becomes stale. Examples of these cases include node failure or nodes leaving the area covered by the system.

Value timeouts limit the time for which any index entry (k, v) of a given value v will be stored in an index cache. By receiving a response from the origin node of (k, v) , the corresponding value timeout will be reset. Let $a_{(k,v)}$ be the time elapsed since (k, v) has been extracted from a response message generated by its origin node. We define the age a_v of value v as $a_v = \min_k (a_{(k,v)})$, i.e., the time elapsed since the most recent response message of this kind was received. When at a node holds $a_v > T$ for the

given timeout value T , all pairs (k, v) are removed from its index cache. PDI implements only one timeout per value v rather than an individual timeout for each index entry (k, v) . This is because in most applications the fact that one index entry (k, v) for a given v expires indicates a substantial change of the value. Subsequently, all other index entries (k', v) are likely to be influenced. For example, in a file sharing system, a pair (keyword_i, URI) is removed when the file specified by URI is withdrawn from the system. Thus, all other pairs (keyword_j, URI) also become stale. Note that depending on the application the concept of value timeouts can be easily extended to individual timeout durations T_v for each value v . Such duration may be included in a response message generated by the origin node. For ease of exposition, we assume in the remainder of this paper a global timeout value T for all values in the system.

To determine the current age of a value, an age field is included in the response message for each value. This age field is set to zero in each response from the origin node. When receiving a response message, a node n extracts the age of each value and calculates the supply time s_v . That is the time at which a response for this value was generated by the origin node. Assume that the response message contains age a_v , then s_v is determined by $s_v = c_n - a_v$, where c_n denotes the local time of node n . s_v is stored in the index cache together with v . Note that v might already be present in the index cache with supply time s'_v . The copy the index cache might result from a more recent response by the origin node, i.e., $s_v < s'_v$. Thus, in order to relate the age of a value to the most current response from the origin node, the supply time is updated only if $s_v > s'_v$. When a node generates a response for a cached index entry (k, v) , it sets the age field for each value v to $a_v = c_n - s_v$. Note that only time differences are transmitted in PDI messages, eliminating the need for synchronizing clocks of all participating nodes.

3.3 Lazy Invalidation Caches for Dealing with Data Modification at Origin Node

Additional to the scenarios described above, a node produces stale index entries by modifying information. That is the case when an index entry is removed from the local index. In a worst-case scenario, a node suddenly leaves the system and all index entries supplied by the node expire at the same time. One way to handle

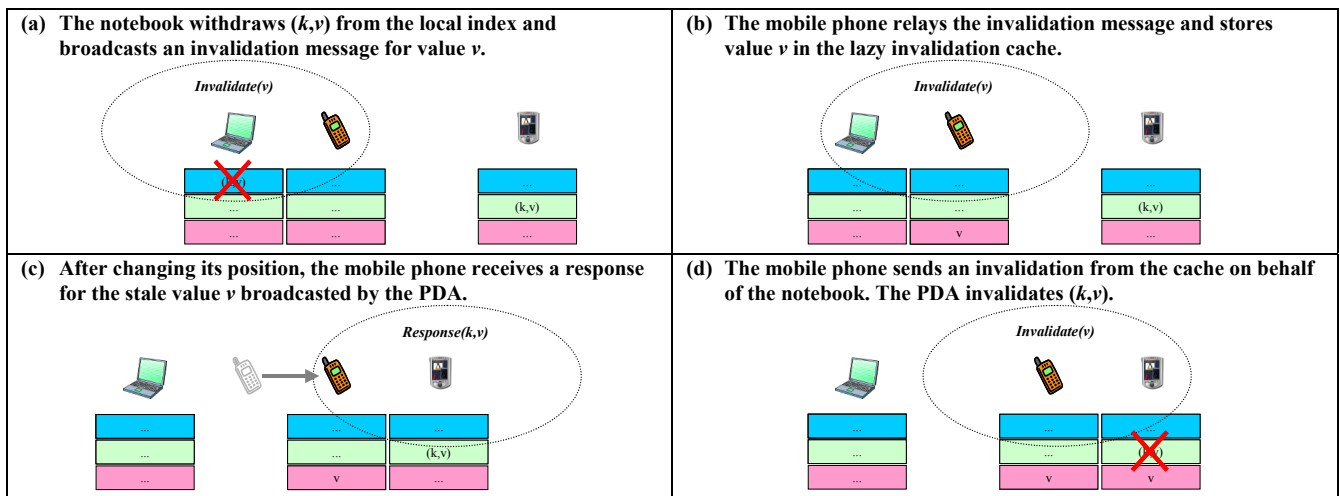


Figure 2. Epidemic dissemination of invalidation messages using lazy invalidation caches

such modification of information is to wait until the timeouts of the values in the stale index entries elapse. Depending on the application and the timeout value T , this straightforward solution may cause severe inconsistency, especially if T is large. A more effective way to handle information modification in distributed applications constitutes the explicit invalidation by control messages.

As basic idea of the explicit invalidation mechanism, a node removes all index entries (k, v) from the index cache when it receives an invalidation message for value v . Flooding is a straightforward way to propagate invalidation messages. To flood an invalidation message for value v , the node removing an index entry sends the message to the limited broadcast address. All mobile nodes that receive the message will relay it exactly once, so that the message is propagated to each node that is connected to the initial node via one or more hops. Unfortunately, in mobile systems even a multi-hop connection between two nodes frequently does not exist. Subsequently, stale index entries are still contained in the index caches of nodes not reached by the invalidation message. Note that the index entries will be redistributed in the system due to the epidemic dissemination. We will show in Section 4.3 that even repeated flooding of invalidation messages does not significantly reduce the number of hits for index entries.

This observation is consistent with [2], which reports that deleted database items “resurrect” in a replicated database environment due to epidemic data dissemination. In [2], a solution is proposed that uses a special message to testify the deletion of an item, referred to as death certificate. Death certificates are actively disseminated along with ordinary data and deleted after a certain time. In contrast, we propose a mostly passive (or “lazy”) approach for the epidemic propagation of invalidation messages, which is illustrated in Figure 2. For the initial propagation of an invalidation message by the origin node, we rely on flooding as described above (Figure 2a). Each node maintains a data structure called lazy invalidation cache, which is drawn as a third box below the mobile devices in Figure 2. When a node receives an invalidation message for a value v it does not only relay it, but stores v in the invalidation cache (Figure 2b). Note that an entry for v is stored in the invalidation cache, regardless if the node stores any index entry (k, v) for v in the index cache. Thus, every node will contribute to the propagation of invalidation messages, so that distribution of information and invalidation messages is separated. To enable the epidemic propagation of the invalidation message, a node scans the invalidation cache for all values contained in an overheard response message (Figure 2c). If a value v' is found, the node will generate an invalidation message for v' itself, because the hit in the invalidation cache indicates that the index cache of a nearby node contains a stale entry (Figure 2d). The invalidation message is not flooded through the complete network, but only with a certain scope similar to forwarding query and response messages as described in Section 2. A node that receives a cached invalidation message will store the included value v in the invalidation cache, and remove all index entries (k, v) from the index cache. Additionally, the node checks whether it has recently received hits for v in response to an own query, which must also be invalidated and may not be passed to the application using PDI.

As the index cache size, the invalidation cache size is limited to a fixed number of values and LRU replacement is employed. In

Section 4.3, we show that setting the invalidation cache size to a fraction below 20% of the index cache size achieves sufficient reduction of false hits assuming a reasonable rate of data modification. Note that LRU replacement does neither guarantee that an invalidation cache entry is kept until all stale index entries are invalidated, nor that it is removed after a certain time, inhibiting a node indefinitely from restoring a value it has once invalidated. Increasing the invalidation cache size solves the first problem, though, doing so amplifies the second problem. To avoid this tradeoff, storing the supply time of invalidation messages similar to the supply time of values as described by Section 3.2 yields an efficient mechanism to decide whether a result for a value is more recent than an invalidation message.

4. Performance Studies

4.1 Simulation Environment

To evaluate the performance of the introduced consistency mechanisms for PDI, we conduct simulation experiments using the network simulator ns-2 [3]. We developed an ns-2 application implementing the basic concepts of PDI as described in Section 2 as well as the two consistency mechanisms described in Section 3. An instance of the PDI application is attached to each simulated mobile device, using the UDP/IP protocol stack and a MAC layer according to the IEEE 802.11 standard for wireless communication [7]. All MAC layer parameters were configured to provide a radio coverage with a radius of 115m. We assume that N_{Nodes} mobile nodes move in an area of 1000 m \times 1000 m according to the random waypoint mobility model [1]. Maximum node speed is 1.5 m/s and a pause time between two movement epochs is 50 s. The random waypoint model is commonly used to mimic the movement of pedestrians. Note that these assumptions clearly constitute conservative assumptions for disseminating information by exploiting mobility.

We considered four different sizes for PDI index caches 32, 128, 512 and 2048 entries. In all experiments, we set the TTL for selective forwarding $TTL_{query} = 4$ hops, and for forwarding of invalidation messages $TTL_{inv} = 2$ hops. We assume that each node contributes $N_{Values} = 16$ values. Each value v_i , $0 \leq i < N_{Nodes} \cdot N_{Values}$, matches a randomly chosen set $K_i \subset \mathbb{N}_0$ of keys. To determine K_i for each value v_i , we adopt the notion of selection power [17], which assumes that the probability that key j matches a given value is $f(j)$. Experiments have shown that $f(j) = (1/\lambda)e^{-j/\lambda}$ for $\lambda \approx 100$ provides a sufficient model of selection power [17]. Following [16], [18], we assume a Zipf-like distribution with $\alpha = 0.9$ for the probability $g(j)$ that a given query is for key j (i.e., $g(j) \approx j^{-0.9}$) and a correlation between query popularity and query selection power, i.e., the most popular key 1 with popularity $g(1)$ has the highest selection power $f(1)$ [17]. Each node issues queries for a randomly chosen key in exponentially distributed intervals with mean 120 s. For simplicity, we restrict queries to single keys, which are chosen randomly according to the popularity distribution $g(j)$.

To cover both spontaneous departures of nodes and data modifications, we assume arrivals and departures of nodes to be exponentially distributed. We considered a ratio between the arrival and the departure rate, so that approximately 30% of all nodes will depart during the considered simulation time. New nodes enter the system with empty index and invalidation caches and contribute values that are not used by any node already in the system. Leaving nodes do not send invalidation messages for

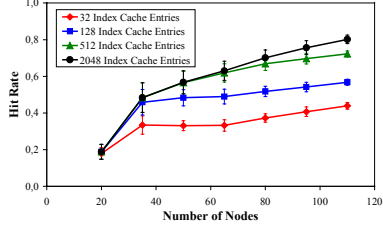


Figure 3. No invalidation: system size vs. hit rate

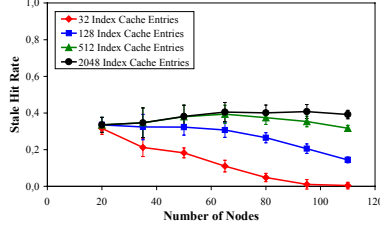


Figure 4. No invalidation: system size vs. stale hit rate

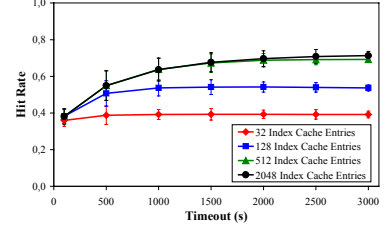


Figure 5. Value timeouts: duration vs. hit rate

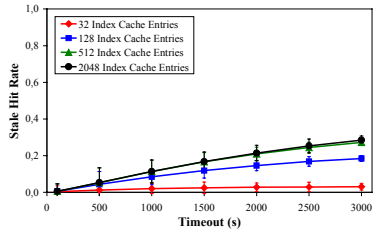


Figure 6. Value timeouts: duration vs. stale hit rate

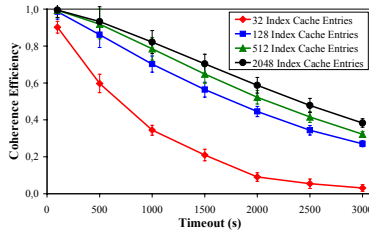


Figure 7. Value timeouts: duration vs. coherence efficiency

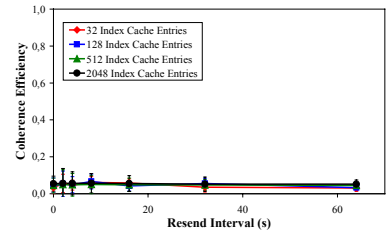


Figure 8. Flooded invalidations: resend-time vs. coherence efficiency

supplied data, i.e., all index entries supplied by a departing node expire. Additionally, each value expires exactly once in a simulation according to a uniform distribution.

We choose performance measures to evaluate the accuracy and the coherence of the results delivered by PDI and the impact of the introduced invalidation mechanisms. Accuracy is measured by the hit rate HR , i.e., $HR = H_F / K_F$ for H_F denoting the number of up-to-date hits and K_F the total number of all up-to-date matching values currently in the system. Note that hit rate can be compared to the information retrieval measure recall. Coherence is measured by the stale hit rate SHR , i.e., $SHR = H_s / (H_s + H_F)$, where H_s denotes the number of stale hits returned on a query. Note that stale hit rate is related to the information retrieval measure precision by $precision = 1 - SHR$. As last measure, reduction of stale hits is measured by the coherence efficiency EF , i.e., $EF = 1 - (H_s / \hat{H}_s)$, where \hat{H}_s denotes the number of hits for stale index entries without employing an invalidation mechanism.

We conduct transient simulations starting with initially empty caches. For each run, the total simulation time was 2 hours. To avoid inaccuracy due to initial warm-ups, we reset all statistic counters after a warm-up period of 10 min. simulation time. For each point in all performance curves, we performed 100 independent simulation runs and calculated corresponding performance measures at the end of the simulation. In all curves 99% confidence intervals determined by independent replicates are included.

4.2 Performance of PDI without Invalidations

In a first experiment, we investigate the coherence of index caches maintained by PDI without the invalidation mechanisms presented in Section 3. These performance curves are shown in Figures 3 and 4. Figure 3 plots hit rates as a function of system size for different sizes of the index cache. The results reveal that hit rate increases with growing system size because an increasing number of nodes increase the dissemination of information. Furthermore, the increase slows down with increasing number of nodes because the total number of values in the overall system

increases. In these cases, the hit rate is clearly limited by the overall index cache size. Note that increasing the index cache size from 512 to 2048 entries, i.e., by factor 4, only increases the hit rate about 10% for systems comprising of many nodes.

Stale hit rates as a function of system size is plotted in Figure 4. We find that without invalidation the stale hit rate is at most 0.4. For smaller index cache sizes, the stale hit rate decreases with system size. Jointly considering Figures 3 and 4 reveals that for increasing system sizes the stale hit rate drops rapidly at the point when the growth of the hit rate slows down. Looking closer at the index caches in these scenarios, we find that the cache content is highly variable. Therefore, stale index entries are removed early from the caches. We conclude from Figure 3 that large caches yield a high amount of stale hits when no invalidation mechanism is used. In contrary, small index caches naturally reduce stale hits, while they fail to provide high hit rates. This evidently illustrates the need for invalidation mechanisms in order to achieve both high hit rates and low stale hit rates.

4.3 Impact of Configurable Value Timeouts

For the following experiment, we fix system size to 80 nodes and investigate the performance of basic PDI extended by value timeouts as implicit invalidation mechanism. Figure 5 plots hit rates versus timeout durations for cache entries. As value timeouts invalidate both stale and up-to-date index entries, the hit rate increases with increasing timeout duration. Thus, invalidations occur less frequently. Unfortunately, as Figure 6 reveals, the stale hit rate increases, too. However, comparing Figures 5 and 6 illustrates that the stale hit rate grows almost linear with an increasing timeout duration while hit rate grows in a log-like fashion. Based on this observation, we choose low timeout duration in order to limit the decrease in hit rate. For example, given an index cache with 2048 entries, a timeout of 1000 seconds decreases the hit rate by about 0.07, while the stale hit rate is decreased from 0.5 to 0.11. That is about 75% improvement compared to the corresponding scenario without invalidation as shown in Figures 3 and 4, respectively. Note that the optimal timeout duration clearly depends on the rate of modification of the information as well as on the arrival and

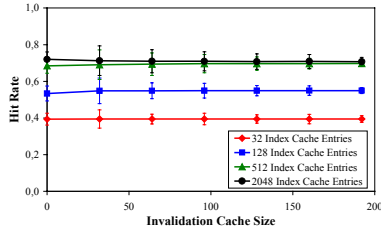


Figure 9. Invalidation caches: cache size vs. hit rate

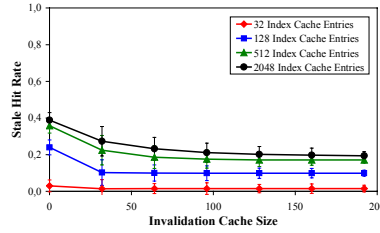


Figure 10. Invalidation caches: cache size vs. stale hit rate

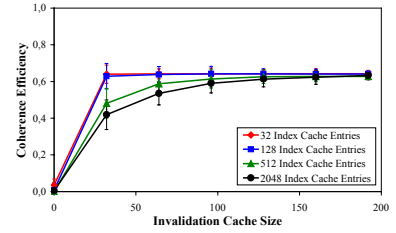


Figure 11. Invalidation caches: cache size vs. coherence efficiency

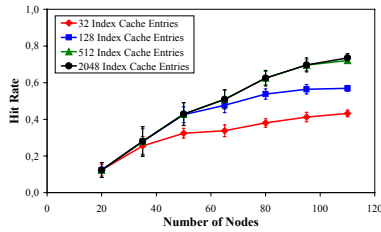


Figure 12. Integrated approach: system size vs. hit rate

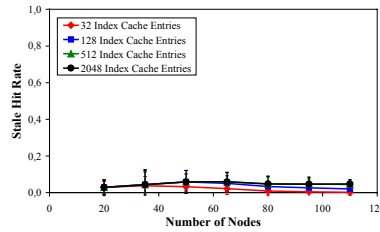


Figure 13. Integrated approach: system size vs. stale hit rate

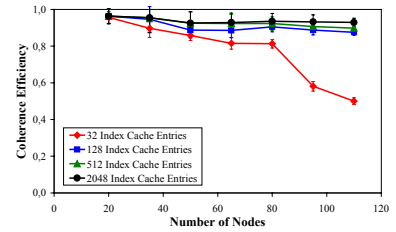


Figure 14. Integrated approach: system size vs. coherence efficiency

departure rates for nodes. Thus, a timeout of 1000s may not be the best choice for all application scenarios. To gain further insight into the behavior of value timeouts, Figure 7 plots the coherence efficiency versus timeout duration. We find that the coherence efficiency rapidly drops with increasing timeout duration due to more infrequently occurring invalidations. Surprisingly, we find that value timeouts are less efficient for small cache sizes than for large ones. Again, this is a confirmation that small caches naturally reduce hits for stale index entries by frequent replacements, shortening the room for improvements by a timeout-based invalidation mechanism. We conclude from Figures 5 to 7 that value timeouts provide an efficient mechanism for implementing implicit invalidation, especially for large index caches.

4.4 Impact of Lazy Invalidation Caches

In the following experiment, we investigate the performance of lazy invalidation caches as explicit invalidation mechanisms. As above, the system size is kept fixed to 80 nodes. Figure 8 plots the coherence efficiency achieved by pure periodical flooding of invalidation messages without using invalidation caches versus re-resent time, i.e., the interval between two successive transmissions of an invalidation message. We find that independent of the re-resent time the coherence efficiency is below 0.1 regardless of index cache size. Note that due to weak connectivity, flooding of invalidation messages cannot eliminate all stale index entries. Subsequently, the epidemic dissemination of remaining index entries leads to a redistribution of stale values. Thus, as illustrated by Figure 8, periodic flooding of invalidation messages fails to implement explicit invalidations.

Subsequently, we investigate performance of the epidemic dissemination of invalidation messages using invalidation caches. In various experiments, we investigated the sensitivity of all performance measures to the TTL_{inv} . Due to space limitations, we do not show these performance results. We found that performance does not significantly increase for $TTL_{inv} > 2$ hops. Thus, we set $TTL_{inv} = 2$ hops in all subsequent experiments. Figure 9 illustrates that lazy invalidation caches regardless of their size do not affect the hit rate. This is because opposed to

value timeouts explicit invalidation messages only invalidate stale index entries.

However, lazy invalidation caches significantly reduce the stale hit rate, especially for large index cache sizes, as shown in Figure 10. We find that for large caches the stale hit rate is reduced by more than 50% compared to Figure 4. Note that if the invalidation cache size increases beyond 20% of the index cache size, no significant further reduction of the stale hit rate can be achieved. For practical applications, this means that the invalidation caches can be small compared to index caches. This observation is confirmed by the results for the coherence efficiency shown in Figure 11. Opposed to the coherence efficiency of value timeouts, the coherence efficiency of lazy invalidation caches is best for small cache sizes. Compared to Figure 7, we find that the coherence efficiency of lazy invalidation caches is smaller than for value timeouts when using large index caches. To understand this observation, recall that nodes leaving the system do not explicitly invalidate all index entries they have supplied, a worst-case scenario for each explicit invalidation mechanism. We conclude from Figures 9 to 11 that lazy invalidation caches efficiently implement explicit invalidation.

4.5 Performance of the Integrated Invalidation Approach

In a last experiment, we investigate the performance of an integrated approach combining both value timeouts and lazy invalidation caches to take into account both weak connectivity and modification of information. We fix the duration of the value timeout to 1000s and the invalidation cache size to 128 entries. Figure 12 plots hit rates versus system sizes. We find that hit rate is reduced mostly for small systems due to invalidations for up-to-date index entries by value timeouts. This leads to a decrease of at most 20%. The performance of index cache sizes of both 512 and 2048 is equal because a large cache cannot benefit from long-term correlations between requests due to the short timeout. For growing system sizes, the hit rate converges towards results without an invalidation mechanism as already shown in Figure 3.

As compensation of the reduction of hit rate, the stale hit rate is significantly reduced compared to a system without invalidation. As shown in Figures 13 and 14, the stale hit rate is highest and coherence efficiency is worst for medium system size and sufficient large index cache sizes. The reason is that a fixed cache size of 128 entries is somewhat to large for small systems, while for large systems the natural limit of stale index entries illustrated in Section 4.2 increases the coherence efficiency. Again, the coherence efficiency drops rapidly for small index cache sizes due to natural reduction of stale hits. We conclude from Figures 12 to 14 that the integrated approach comprising of the introduced implicit and explicit invalidation mechanisms can effectively handle both spontaneous node departures and modification of information. In fact, for large index caches, the stale hit rate can be reduced by more than 85%.

Conclusion

We introduced a distributed lookup service for mobile applications denoted as Passive Distributed Indexing (PDI) and presented invalidation mechanisms for reducing inconsistency in PDI index caches. PDI stores entries in form of (key, value) pairs in index caches located in each mobile device. For implementing implicit invalidation of index entries, value timeouts were proposed. The epidemic distribution of invalidation messages based on lazy invalidation caches implemented explicit invalidation. Presented simulation studies using ns-2 showed that value timeouts reduce inconsistent results due to intermitted connectivity or node failure by more than 75%. The employment of lazy invalidation caches reduces the fraction of stale index results by more than 50%. The integration of both invalidation mechanisms achieves a coherence efficiency of more 95% for the considered scenario.

Recall that in the simulation studies, we considered as a conservative assumption a mobility model for mimicking the movement of pedestrians. The spread of information by PDI is considerably faster for higher mobility. Thus, PDI is likely to achieve considerable higher hit rates while providing a comparable performance in terms of coherence efficiency for such scenarios. Furthermore, we conjecture that optimizing the duration for value timeouts and the size for lazy invalidation caches will further improve hit rate.

References

- [1] J. Broch, D. Maltz, D. Johnson, Y.-C. Hu, and J. Jetcheva, A Performance Comparison of Multi-Hop Wireless Ad Hoc Network Routing Protocols, *Proc. 6th ACM/IEEE MobiCom 98*, Dallas, TX, 85-97, 1998.
- [2] A. Demers, D. Greene, C. Hauser, W. Irish, J. Larson, S. Shenker, H. Sturgis, D. Swinehart, and D. Terry, Epidemic Algorithms for Replicated Database Maintenance, *Proc. 6th Symp. on Principles of Distributed Computing (PODC 1987)*, Vancouver, Canada, 1-12, 1987.
- [3] K. Fall and K. Varadhan (editors), The ns-2 manual, *Technical Report*, The VINT Project, UC Berkeley, LBL, and Xerox PARC, 2003.
- [4] S. Goel, M. Singh, D. Xu, and B. Li, Efficient Peer-to-Peer Data Dissemination in Mobile Ad-Hoc Networks, *Proc. Int. Workshop on Ad Hoc Networking (IWAHN 2002)*, Vancouver, BC, 2002.
- [5] W. Heinzelman, J. Kulik, and H. Balakrishnan, Adaptive Protocols for Information Dissemination in Wireless Sensor Networks, *Proc. 5th ACM/IEEE MobiCom 99*, Seattle, WA, 174-185, 1999.
- [6] K. Hanna, B. Levine, and R. Manmatha, Mobile Distributed Information Retrieval For Highly-Partitioned Networks, *Proc. 11th IEEE Int. Conf. on Network Protocols (ICPN 2003)*, Atlanta, GA, 2003.
- [7] IEEE Computer Society LAN MAN Standards Committee, Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications, *IEEE Standard 802.11-1997*, New York, NY, 1997.
- [8] A. Khelil, C. Becker, J. Tian, and K. Rothermel, An Epidemic Model for Information Diffusion in MANETs, *Proc. 5th ACM Int. Workshop on Modeling, Analysis and Simulation of Wireless and Mobile Systems (MSWiM 2002)*, Atlanta, Georgia, 2002.
- [9] C. Lindemann and O. Waldhorst, A Distributed Search Service for Peer-to-Peer File Sharing in Mobile Applications, *Proc. 2nd IEEE Conf. on Peer-to-Peer Computing (P2P 2002)*, Linköping, Sweden, 71-83, 2002.
- [10] Internet Engineering Task Force Working Group Mobile Ad hoc Networks (MANET). <http://www.ietf.org/html.charters/manet-charter.html>.
- [11] P. Mockapetris, Domain Names - Concepts and Facilities, *IETF Request for Comments 1034*, 1987.
- [12] M. Papadopouli and H. Schulzrinne, Effects of Power Conservation, Wireless Coverage and Cooperation on Data Dissemination among Mobile Devices, *Proc. 2nd ACM MobiHoc 2001*, Long Beach, NY, 117-127, 2001.
- [13] M. Papadopouli and H. Schulzrinne, Performance of Data Dissemination and Message Relaying in Mobile Ad Hoc Networks, *Technical Report CUCS-004-02*, Columbia University, 2003. Under submission.
- [14] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker, A Scalable Content-Addressable Network, *Proc. ACM SIGCOMM 2001*, San Diego, CA., 149-160, 2001.
- [15] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan, Chord: A Scalable Peer-to-Peer Lookup Service for Internet Applications, *Proc. ACM SIGCOMM 2001*, San Diego, CA, 149-160, 2001.
- [16] K. Sripanidkulchai, The Popularity of Gnutella Queries and its Implications on Scalability. *Proc. O'Reilly Peer-to-Peer and Web Services Conf.*, 2001
- [17] B. Yang and H. Garcia-Molina, Comparing Hybrid Peer-to-Peer Systems, *Proc. 27th Int. Conf. on Very Large Data Bases*, Rome, Italy, 561-570, 2001.
- [18] Y. Xie and D. O'Hallaron, Locality in Search Engine Queries and Its Implications for Caching, *Proc. IEEE INFOCOM 2002*, New York, NJ, 2002.