# Update Propagation Through Replica Chain in Decentralized and Unstructured P2P Systems

Zhijun Wang, Sajal K. Das, Mohan Kumar and Huaping Shen Center for Research in Wireless Mobility and Networking (CReWMaN) Department of Computer Science and Engineering The University of Texas at Arlington, Arlington, TX 76019, USA Email: {zwang, das, kumar, hpshen}@cse.uta.edu

# Abstract

In this paper, we propose a novel algorithm, called update propagation through replica chain (UPTReC), to maintain file consistency in decentralized and unstructured peer-to-peer (P2P) systems. In UPTReC, each file has a logical replica chain composed of all replica peers (RPs) which are defined as peers that have replicas of the file. Each RP acquires partial knowledge of the bi-directional chain by keeping a list of information about k nearest RPs, called probe peers, in each direction. When an RP initiates an update, it pushes the update to all possible online (active) RPs through the replica chain. A reconnected RP pulls an online RP to synchronize the replica status and the information of the probe peers. An analytical model is derived to evaluate the performance of the UPTReC algorithm. The analytical results provide a better understanding of the system in choosing the system parameters for probabilistically guaranteed file consistency with minimum overheads. Simulation experiments are conducted to compare the performance with an existing update propagation algorithm based on the rumor spreading scheme. The experimental results show that the UPTReC can significantly reduce (up to 70%) overhead messages and also achieve smaller stale query ratio for files prone to frequent updates.

# I. INTRODUCTION

The peer-to-peer (P2P) systems are self-organizing distributed systems, in which all participating peers cooperatively provide and receive services from each other. P2P systems are rapidly growing due to such desirable characteristics as scalability, availability, anonymity and authentication.

The P2P systems can be categorized into structured and unstructured. In a structured P2P system, the topology is tightly controlled and the files are well deployed [14] [16]. On the other hand, an unstructured P2P system has no central control of its topology and file placement [1][3] [5]-[7] [9] [12]. Chord [16] is an example of a structured P2P system whereas Gnutella [1] is an example of a decentralized and unstructured P2P system. In this paper, we focus on the decentralized and unstructured P2P systems, in which each peer maintains information about its neighboring peers. Some files may be heavily replicated in the system to improve the file availability and system fault-tolerance. To acquire a file, a peer searches the file through its neighboring peers. In order to make such P2P systems scalable and efficient, significant efforts have been made on the development of search and replication algorithms [5] [7][9][12]. In these algorithms, the locations of replicas for a file are well deployed based on partial knowledge of the system to minimize the search cost and balance the network load. Furthermore, these algorithms assume that the files are rather static and updates occur very infrequently. Indeed, the impact of the file update has not received much attention.

However, in many application domains, such as trust management [3], bulletin-board systems and distributed web cache [10], the updates occur frequently. Let us look at a simple file update example by considering peers sharing software. Assume a peer initiates a software in the system, which can be replicated by other peers to enhance its availability and minimize search costs. Any peer which has the replica of the software can update it, as a result of debugging or adding more advanced features. After a peer modifies the software, its old version replicas are no longer valid, hence if the update is not properly propagated to its replicas, the incoming queries from these replicas are not valid. Without update propagation, even if the replicas are invalidated properly, the well placed replicas no longer exist, thus resulting in a large search cost or even unavailable to find the file for incoming file queries. Therefore, effective propagation of update to all replica peers (RPs), which are defined as peers that have the replicated files, is critical to maintaining balanced network load, enhanced file availability, and reduced access latency. Datta, et al. [8] proposed a hybrid push/pull update propagation algorithm based on the rumor spreading algorithm for decentralized and unstructured P2P systems. This algorithm is the first attempt to focus on the effective propagation of updates in such P2P systems. However, the overhead messages of update propagation are significant, implying significant resource (bandwidth) consumption for frequently updating files in large P2P systems. Moreover, it is difficult to deal with RPs with dynamic IP addresses.

In this paper, we propose a novel algorithm, called Update **P**ropagation Through **Re**plica Chain (UPTReC), to maintain file consistency in decentralized and unstructured P2P systems.

UPTReC provides a probabilistically guaranteed file consistency. In this algorithm, each file has a logical replica chain composed of all RPs. Each RP has partial knowledge of the bi-directional chain by keeping information (i.e., identity (ID) and IP address) of a list of k nearest RPs (called *probe* peers) in each direction. The replica chain can be naturally built and easily maintained during the file replica process. When an RP updates a file, it pushes the update to its online (i.e., active) probe peers in each direction. The farthest online probe peer in turn forwards the update to its online probe peers along the direction. This process recursively propagates the update to all possible online RPs. When an offline (i.e., inactive) RP gets reconnected, it pulls an online probe peer to synchronize the file status and the probe peers' information. If the RP's IP address is changed, the new IP address is pushed to all possible online probe peers which in turn update the maintained information of the reconnected RP.

An analytical model is derived for the proposed UPTReC algorithm. The analytical results provide a better understanding of the system in choosing the system parameters for probabilistically guaranteed file consistency with minimum overheads. The analytical results also show that UPTReC is efficient for RPs with dynamic IP addresses. The simulation results show that UPTReC significantly reduces (up to 70%) overhead messages to propagate updates in comparison with the rumor spreading based algorithm [8].

The rest of the paper is organized as follows. Section II gives an overview of the related work. A detailed description of UP-TReC is given in Section III. An analytical model is derived in Section IV. Section V presents performance comparisons of UPTReC with an existing propagation algorithm. The conclusions are drawn in Section VI.

#### II. RELATED WORK

The problem of searching and replicating files in P2P systems has received much attention. However, most P2P systems consider files to be static and do not address the issue of file updates and file consistency maintenance. In this section, we present an overview of algorithms that focus on file consistency maintenance in decentralized and unstructured P2P systems.

Datta et al. [8] proposed a hybrid push/pull update propagation algorithm based on the rumor spreading algorithm for highly unreliable and unstructured P2P systems, such as Gnutella [1] and P-Grid [2]. The algorithm provides probabilistically guarantees rather than strict consistency. Here, each RP maintains a subset of all RPs as its responsible peers. When an RP initiates an update, the update is pushed to its responsible peers, which in turn propagate the update to their responsible peers with some probabilities. This process continues until all possible online peers get the update. When a peer gets reconnected, it queries multiple responsible peers to synchronize itself with the peer having the most recent update.

The algorithm in [8] is the first attempt to focus on the effective propagation of updates to RPs for decentralized and unstructured P2P systems. However, the overhead messages due to push updates are significant. Moreover, the maintenance of the subset of responsible peers is not easy, especially for RPs with dynamic IP addresses. There is no discussion on the responsible subset maintenance in [8].

An invalidation report based on push and pull (PAP) algorithm is developed by Lan, et. al [11]. In PAP, each file has a master peer, only the master peer can update the file. An estimated Time-To-Expire (TTE) and the master peer information are associated with each replica. When a file is updated, its invalidation report is broadcast to the network. Any online peers that have replicas of the file invalidate the replicas. Once the TTE of a file expires, the file must be pulled from the master peer if it is accessed. Only the master peer updating the file is a strong constraint in P2P systems. Moreover, the master peer may change its IP address and go offline, thus resulting in a small probability of an RP successfully pulling a master peer.

# III. UPDATE PROPAGATION THROUGH REPLICA CHAIN (UPTREC)

The main motivation behind UPTReC is to minimize the overhead messages for propagating updates to RPs in decentralized and unstructured P2P systems. The detailed description of the UPTReC algorithm is given in the following subsections.

#### A. System Model and Assumptions

We consider a decentralized and unstructured P2P system, such as Gnutella where all peers are equal and no peer has a global view of the system. The system model and assumptions are summarized as follows:

- 1) No strong file consistency is required, but a probabilistically guaranteed file consistency is required.
- 2) The write-write conflict is ignored.
- 3) All peers frequently join and leave the system.
- 4) An online peer that gets an update has the ability to finish its push process.
- 5) An online peer can communicate with any other online peer if it knows the IP address of that peer.
- 6) The physical connectivity and system topology are ignored.
- 7) Each RP has an ID and an IP address, the ID is fixed but the IP address may be changed for each reconnection.
- 8) Each file is associated with a version and generation time used for synchronization.

In UPTReC, if two RPs update the file and push it through the chain at the same time, an RP in the chain can detect a write-write conflict when it receives two updated files of the same version generated by two different RPs. In this case, the RP that detected the conflict can send the conflict information back to the two update generating RPs, which in turn solve the conflict through communication with each other, then the latest updated file is pushed through the chain again. Due to very lower write-write conflict rate [13] in P2P systems, we make assumption (2).

The probability of online peer to successfully finish its push process is usually over 0.95 [8]. If the probability is low for a

system, the assumption (3) above can be remedied by using a reliable push process. In a reliable push process, the push process of RP a does not stop after it propagates the update to an online RP b, which in turn forwards the update to other RPs. RP a must wait for the confirmation from RP b indicating the update has been successfully propagated. If the confirmation is not received within a certain period, RP a probes RP b again; and if RP b is offline, RP a contacts with another RP to continue the push process. The reliable push process incurs some additional overhead messages for confirmation. We make assumption (3) to simplify our algorithm analysis.

#### B. Push Update Through Replica Chain

Figure 1 (a) shows a logical replica chain for a file with N replicas. Each RP is a node on the chain and has a unique ID associated with it. Each node <sup>1</sup> maintains information (i.e., ID and IP address) about k (typically k is tens) nearest nodes in each (left and right) direction of the chain <sup>2</sup>. These 2k nodes are called as *probe* nodes (peers). Two nodes are said to have h-hop distance if there are h-1 nodes between them. For example, node i and node i+k have k-hop distance.



Fig. 1. (a) Logical Replica Chain; (b) Update propagation of node *i*. Node  $i \pm m (0 < m \le k)$  is the  $m^{th}$  probe node in right (left) direction.

Figure 1 (b) shows the update propagation process of node *i*. When node *i* initiates an update (node *i* is called the update initiating node), the update is pushed symmetrically along both left and right directions of the chain. Now let us look at the process of node i pushing the update to node N (right side of the chain). Node i has information of k probe nodes in right direction (from node i+1, called the  $1^{st}$  probe node, to node i+k, called the  $k^{th}$  probe node). To push an update, node isends a probe message to each of the probe nodes in this direction (i.e., node i+k, ..., i+1). The farthest online probe node (here node i+k-1) is chosen to be the update relay node, which will further propagate the update through the chain along the direction. All other online probe nodes of i, such as node i+k-2, will receive but do not propagate the update. After node idetermines its update relay node i+k-1, it first sends the update to that node with the relay flag bit set as 1 and then sends

the update to all other online probe nodes with the relay flag bit set as 0. When an online probe node receives the update, it first checks the update relay flag bit. If the bit is 0, it only needs to receive the update. Otherwise, it needs to propagate the update through the chain along the direction. The process of the update propagation is similar to node *i* except not to send the probe messages to its probe nodes which are also the probe nodes of *i*. Because all these nodes are probed by *i* and they should be offline. As shown in Fig. 1(b), when node i+k-1 gets the update, it finds that the update relay flag bit is 1, and hence it immediately sends the probe messages to its probe nodes in the right hand side which are not the probe nodes of *i*, i.e., nodes  $i+2k-1, \dots, i+k+1$ . The update propagation process is repeatedly executed through the replica chain. If all k probe nodes of an update relay node are offline, the propagation process is stopped and the update can not be propagated in this direction. The same process is executed for node *i* to propagate the update to node 1.

#### C. Pull After Online

During the offline period of an RP, it may miss some updates of the file and/or some information on the chain changes. Hence, when an offline RP gets reconnected, it needs to pull some online RPs to synchronize the status of the file and its probe nodes. An RP can probe an online probe node from its nearest probe node to farthest one in each direction. Whenever an online RP is probed in one direction, the file and the information of its probe nodes are synchronized. If its IP address is not changed, the pull process in this direction is finished. The same process is executed in the other direction. If the IP address of the reconnected RP is changed, it needs to send its ID and new IP address to all its probe nodes. Then the pull process is finished. If no online probe node can be pulled (due to probe nodes going offline or changing IP addresses), the reconnected RP needs to connect its probe nodes through flooding search to synchronize the status of the file and the information of its probe nodes if its IP address is changed.

### D. Chain Construction and Maintenance

We discuss how to construct and maintain the replica chain in this subsection. After a peer initiates a file in the system, the file can be searched, fetched and replicated by other peers. Each replica is copied from one of the other replicas. If each RP maintains the information of all RPs which fetched a file from it, then a replica tree is naturally constructed. Figure 2 (a) shows a replica tree composed of 5 RPs as the root node at RP 1. If all RPs are always online, any update from any RP can be successfully propagated to any other RPs. For example, when RP 3 initiates an update, it sends the update to RPs 1, 4 and 5. Each RP in turn updates its replica and then relays the update to all its children and parent except the one which sent the update. The update is successfully propagated through all RPs. A new replica tree with RP 3 as the root is shown in Figure 2 (b). However, frequently disconnected peers make such a replica tree ineffective in terms of update delivery. In

<sup>&</sup>lt;sup>1</sup>For simplify, both the node and RP can be used as the RP in the following of the paper.

<sup>&</sup>lt;sup>2</sup>Note that the nodes at or near the head or tail have less than k probe nodes in one direction

order to increase the probability of successfully propagating the update, each RP must maintain the information of multiple RPs along each path. Due to the properties of the general tree, some RPs may maintain the information of a large number of RPs, while some other RPs maintain information of very few RPs. To balance the overhead associated with the file maintained by each RP, a replica chain can be constructed from the replica tree as explained below.



Fig. 2. RPs naturally constructs a replica tree. (a) Root at RP 1; (b) New tree with root at RP 3.

Figure 3 shows the process of constructing a replica chain during the replica process. Figure 3 (a) presents the first four RPs which naturally form a chain. In this case, a new node locates at the head or tail of the chain. When a new peer replicates the file fetched from another RP, the corresponding chain information is also fetched. The information of a new RP is forwarded to all possible RPs which should have the information of the new RP. Figure 3 (a) illustrates the process for an RP, such as RP 4 joining the chain. When RP 4 fetches the file from RP 3, the replica chain including information about RPs 1 and 2 is also fetched. The RP 3 adds RP 4 into the chain, and pushes information about RP 4 to RP 1 and 2. However, if RP 1 for example is offline at that time, it needs to probe either RP 2 or RP 3 to get the latest chain information.



Fig. 3. The process for construction a replica chain. (a) The new replica locates on the head or tail of the chain; (b) The new replica locates at the middle of the chain.

If a new peer joins in the middle of a chain, it needs to push its information to at most k RPs in the chain along the direction opposite to the RP which provides the file. For example, when RP 5 joins the chain by obtaining the chain information from RP 3, RP 5 pushes the information to RP 4.

When RP i removes a replicated file, it sends a message to each of its probe peers to get removed from the replica chain. All online probe peers get the message and in turn remove RP i from the chain. All offline probe peers get this message when

they reconnect. If all probe peers are offline, RP i is not removed from the chain and informs the reconnecting peers when they probe. The process of adding or removing a replica requires up to 2k messages.

#### **IV. PERFORMANCE ANALYSIS**

An analytical model is developed in this section. One critical issue concerning the UPTReC algorithm is to determine the value of k. If k is too small, an update may fail to be propagated through the chain. If k is too large, the overhead cost of chain maintenance is high.

# A. Performance Analysis

Our analytical modeling is based on the assumptions made in Section III-A. Some parameters and measurement metrics are defined in the Table I.

- N: number of RPs in the chain, i.e., the total number of replicas for a file.
- k: number of probe RPs in one direction.
- *P*<sub>on</sub>: probability of an RP to be online.
- $P_{off}$ : probability of an RP to be offline  $(P_{off}=1-P_{on})$ .
- *P*<sub>cIP</sub>: probability of an RP to change the IP address after reconnecting.
- *h*: number of hops an online RP from the update initiating peer.
- T: average period of a peer online and offline cycle.
- $\lambda$ : access rate of a file for the whole system.
- $T_{up}$ : average file update period
- $P_h^s$ :probability of successfully propagating an update to an online RP with *h*-hop distance.
- $P_h^s(m)$ : probability of successfully propagating an update to an online RP with *h*-hop distance while the online RP only counts the contributions of its *m* farthest  $(1 \le m \le k)$  probe peers (i.e., the  $k^{th}$ , ...,  $(k-m+1)^{th}$  probe peers).
- $P_{pull}^{s}(k)$ : probability of a reconnected RP to successfully pull an online RP.
- $C_{flood}$ : average number of messages to find an online probe peer through flooding search.
- $C_{push}(N)$ : maximum number of messages to push an update through an N-node replica chain.
- C<sub>pull</sub>(k): average number of messages in each pull procedure of a reconnected peer.
- *OHQ*: number of overhead messages per query of file consistency maintenance (including overhead of push and pull).
- $P_{stale}(N)$ : stale query probability for a file with N replicas.

In UPTReC, the maximum number of messages to push an update is N, because each RP at most receives one probe message. Thus we have

$$C_{push}(N) \le N \tag{1}$$

When an offline RP rejoins the system, it pulls an online RP from its probe peers in each direction to synchronize the file

status and probe peers' information, the pull process in one direction stops whenever an online RP is pulled. If a probe peer is offline or online but with different IP address from the reconnected RP maintained, it can not be pulled. We use  $P_{fail}=P_{off}+P_{on}P_{cIP}$  to represent the probability that a probe peer can not be pulled by a reconnected peer. Then the probability of a reconnected RP to successfully pull an online RP is

$$P_{null}^{s}(k) = 1 - (P_{fail})^{2k}$$
(2)

If the IP address of the reconnected RP is changed, it needs to contact with all probe peers once, hence 2k probe messages are needed. If no probe peer is probed, it needs to search a probe peer through flooding. So the average number of probe messages for each pull process is:

$$C_{pull}(k) = P_{cIP}[2k + (1 - P_{pull}^{s}(k))C_{flood}] + 2(1 - P_{cIP})[(1 - P_{fail})\sum_{i=1}^{k-1} i(P_{fail})^{i-1} + k(P_{fail})^{k-1})] = P_{cIP}[2k + (1 - P_{pull}^{s}(k))C_{flood}] + 2(1 - P_{cIP})(\frac{1 - P_{fail}^{k}}{1 - P_{fail}})$$
(3)

In Equation 3, the first term is the pull cost for a reconnected RP with changed IP address, and the second term is the pull cost when its IP address is not changed. In the second term, if an online probe peer is pulled in a direction, the pull process is stopped in that direction, and if no online probe peer is pulled, all *k* probe peers are needed to be pulled once. The pull process is symmetrical in both directions.



Fig. 4. Calculation diagram of  $P_h^s(m)$ .

Based on the definitions, we have  $P_h^s = P_h^s(k)$ , and  $P_h^s(m)$  can be recursively calculated. Figure 4 shows the calculation diagram of  $P_h^s(m)$ . Here RP *h* has *h*-hop distance from the update initiating peer.  $P_h^s(m)$  represents the probability for RP *h* to get the update if only its farthest *m* probe peers (i.e., its *k*-th, (*k*-1)-th, ..., (*k*-*m*)-th probe peers) are considered, these probe peers are *h*-*k*, *h*-*k*+1, ..., *h*-*k*+*m*-1 hops distance from the update initiating peer, we call these RPs as RPs *h*-*k*, *h*-*k*+1, ..., and *h*-*k*+*m*-1 as shown in Figure 4. For example,  $P_h^s(2)$  is the probability for RP *h* to get an update if only probe peers *h*-*k* and *h*-*k*+1 are considered to push the update to peer *h*, and all probe peers *h*-*k*+2, ..., *h*-1 are not considered. All these probabilities can be recursively calculated by the following three equations:

If  $h \leq k$  and  $1 \leq m \leq k$ ,

$$P_h^s(m) = 1 \tag{4}$$

If h > k and m = 1,

$$P_h^s(m) = P_{on} P_{h-k}^s(k) \tag{5}$$

If h > k and  $1 < m \leq k$ ,

$$P_{h}^{s}(m) = P_{h}^{s}(m-1) + P_{on}P_{off}^{m-1}P_{h-k+m-1}^{s}(k-m+1)$$
(6)

Equation (4) means that an online RP is a probe peer of the update initiating peer, it can absolutely get the update. Equation (5) indicates that only considering its farthest probe peer h-k, if it is online and successfully receives the update, then RP h can successfully get the update. Equation (6) can be explained by considering the  $m^{th}$  farthest probe peer h-k+m-1, the probability of successfully receiving the update by peer h is the probability of successfully receiving the update through its farthest m-1 probe peers plus the contribution of the  $m^{th}$  farthest probe peer. The  $m^{th}$  probe peer has contributions only if all farthest m-1 probe peers are offline, because if any of these peer is online, the contribution has been counted through that peer. In this case, the probability of successfully getting the update for the  $m^{th}$  farthest probe peer is only through its k-m+1 probe peers (its first m-1 probe peers are offline), i.e.,  $P^s_{h-k+m-1}(k$ -m+1).

The number of overhead message per query of file consistency maintenance is:

$$OHQ = \frac{1}{\lambda} \left(\frac{C_{push}}{T_{up}} + N \frac{C_{pull}(k)}{T}\right)$$
(7)

For a replica chain with N RPs, the maximum number of hops from an update initiating peer to an online RP is N-1. Hence any online RP has a probability larger than  $P_N^s$  to get the update. An offline RP has  $P_{pull}^s(k)$  probability to synchronize with an online RP, hence each online RP has at least  $P_N^s P_{pull}^s$  probability with a valid file. Then the stale query probability is upper bounded by:

$$P_{stale}(N) \le 1 - P_N^s(k) P_{pull}^s(k) \tag{8}$$

The performance of UPTReC is formulated by equations (1) - (8). All these measurements are determined by  $P_{on}$ ,  $P_{cIP}$ , k and N.

#### B. Numerical Results

Some numerical results are shown in this subsection to characterize typical value of k under some probabilistically guaranteed file consistency. The difference between the numerical and simulation results (not presented in the paper) is within 2% in all these cases.

1) Probability of successfully propagating an update through the chain: We study the impact of the number of probe peers (k) on the probability  $(P_h^s)$  of successfully propagating an update to an online RP with h = 10,000 hops. The relationship between  $P_h^s$  and k is shown in Figure 5. When  $P_{on} \ge 20\%$ ,  $P_h^s$  is very close to 1 for  $k \ge 60$ . To achieve  $P_h^s$  close

to 1, k = 40 is enough for  $P_{on} = 30\%$  and k is reduced to 20 for  $P_{on} = 50\%$ . For very small  $P_{on} = 10\%$ , we get k = 110. The results indicate that k = 60 ensures a near to 1 probability to propagate an update through a 10,000-node chain for  $P_{on} \ge 20\%$ . As stated in the previous section, a larger k leads to more overhead messages for the replica chain maintenance. But the overheads per update propagation is independent on k.



Fig. 5.  $P_h^s$  versus k

2) Scalability on the number of replicas: The maximum number of hops of a replica chain increases as the number of RPs increases. In P2P systems, the typical number of replicas for a file varies from tens to thousands. We investigate the scalability of the algorithm on the number of RPs. Figure 6 shows the results of  $P_h^s$  as h increases from 1,000 to 1,000,000. For a system composed of peers with high online probability ( $P_{on} \geq 50\%$ ), a small number of probe peers k = 20 can ensure a larger than 0.95 probability of successful propagation an update to an online RP with 1,000,000-hop distance. For a system with very low online probability RPs, k = 120 makes  $P_h^s > 0.98$  for h = 1,000,000. The probability of successful propagation drops slowly as the number of hops increases. The results indicate that UPTReC algorithm has good scalability in terms of the number of RPs.



Fig. 6.  $P_h^s(k)$  versus h

#### V. PERFORMANCE COMPARISONS

The performance comparisons between UPTReC and the update propagation algorithm based on the rumor spreading algorithm (in short, *Rumor*) proposed in [8] are presented in this section. The overhead messages of file consistency maintenance come from push and pull processes, the major messages of a fast (slow) updating file is from the push (pull) process. We use simulations to study the impact on the performance of update frequency that is not analyzed in Rumor algorithm.

The both algorithms, i.e., UPTReC and Rumor, focus on the efficient update propagation to all online RPs. Note that the update propagation is only through the RPs. Moreover, both algorithms are independent on the file search and replication. Therefore, we simulate only RPs instead of a whole P2P system to focus on the file consistency maintenance cost. The system topology and physical connectivity are ignored.

In the simulations, each RP alternatively leaves and joins the system as a Poisson process. The file update is also assumed to be a Possion process. When an update comes, the update initiating peer is randomly chosen from an online RP. In a real P2P system, a file can be searched and replicated by other peers, and an RP may drop a replica. As stated in the previous section, adding a new RP or removing an RP costs 2k messages to maintain the chain, but the subset maintenance is not discussed in Rumor algorithm [8]. Hence, we ignore the comparison on the costs of the chain and subset maintenance in the simulation by assuming a static chain and subsets. Moreover, all RPs are considered to have static IP addresses, because no method is discussed to deal with dynamic IP address in Rumor algorithm. The chain is randomly built, i.e., each RP has equal probability to appear at any location on the chain. Each RP keeps information of k probe peers in each direction. In the Rumor algorithm, each peer randomly picks up R RPs as its responsible peers. In the 0 push round, the update as well as a replica list are forwarded to its all responsible peers. The replica list records all RPs in which the update has been sent. In the  $t \ (\geq 1)$  push round, a peer has a probability  $P_F(t) = f^t$  to push the update to its any responsible peer that is not on the replica list, where f is a constant between 0 and 1. A RP that receives an update is assumed to have ability to finish its push process. The pull process in both algorithms is similar. In UPTReC, when an online probe peer is probed in a direction, the pull process in this direction is finished. In Rumor, two online probe peers are probed in each pull process.

Let the file have an access rate  $\lambda$  for the whole system, each access randomly fetches the file from an online RP. When an online RP answers a query, if the file is generated in its newest version, a valid query is counted, otherwise a stale query is counted. Due to focus on the efficiency of file consistency maintenance, the parameters  $\lambda$ , T, and  $T_{up}$  are set to unit time.

In our simulation model, when an RP a pushes an update to RP b, it first probes RP b. If RP b is online, the update is forwarded. Thus the total number of update sent out is equal to the number of online RPs which have received the update. This number is almost equal in both algorithms if the stale query ratio is close to each other. We compare the overhead messages

TABLE I Parameter Setup I

N	λ	T	$T_{up}$
10000	1	10000	10000

for probing all RPs rather than the number of update themselves. Of course, the update can be sent out instead of the probe messages. However, if the update is large, this may cause large extra traffic for sending the update to offline RPs.

#### A. Overhead messages for each push process

The number of overhead message in the push process and the stale query ratio are studied under various probabilities of online peers. The probability of successfully propagating an update is determined by the probability of a peer being online and the number of probe (responsible) peers. Based on the analytical results in the previous section, we set  $2kP_{on} = 20$  (or  $RP_{on} = 20$ ) to ensure a low stale hit probability. Thus,  $P_{on} =$ 10% corresponds to k = 100 (R = 200), and  $P_{on} = 50\%$  corresponds to k = 20 (R = 40). The other parameters are set as in Table I. Based on these setups, there are 1 query per RP and 1 update in each RP online and offline cycle (T period) on the average. Two different f values (0.8 and 0.9) are used in the Rumor algorithm to show the relationship between the stale query ratio and the number of overhead messages. The number of overhead messages in Rumor algorithm is determined by the stale query ratio, a larger f or R makes a lower stale query ratio. We set the R value as 2k which is the total number of probe peers kept by an RP in UPTReC. For such R value, high f values are needed to ensure a similar stale query ratio between UPTReC and *Rumor* algorithm, so f is set to 0.8 and 0.9.

Figures 7 and 8 show the number of overhead messages in the push process and the stale query ratio of both algorithms. As shown in these figures, a smaller f reduces the number of overhead messages in the Rumor algorithm, but the stale query ratio is increased. When f drops from 0.9 to 0.8, the number of overhead messages drops about 20%, but the stale query ratio is almost doubled.



Fig. 7. The number of overhead messages for push process versus peer online probability



Fig. 8. Stale query ratio (%) versus peer online probability

The results show that the number of push overhead messages divided by the number of RPs (N) in UPTReC is almost 1, and this value is more than 2.4 in Rumor. The stale query ratio for the UPTReC is less than 1.2% for all ranges of  $P_{on}$  from 10% to 50%. But the stale query ratio for Rumor increases from 1.2% to 4.8% when  $P_{on}$  increases from 10% to 50% with f = 0.8. The stale query ratio can be reduced to less than 2.5% but incurs more than 20% overhead messages if f is set to 0.9. The results indicate that compared with Rumor, UPTReC reduces more than 60% overhead messages to put an update while achieving a smaller stale query ratio.

# B. Overhead messages per query

The number of overhead messages per query in various update frequency is investigated in this case. We measure two performance metrics: the number of overhead messages per query and stale query ratio. The number of overhead messages per query is defined as the total number of consistency maintenance messages which include overhead messages of the push and pull processes divided by the total number of queries in the system. We set two R values (80 and 100) for Rumor algorithm in the simulation to study the effects of R. The other system parameters are set as in Table II.

TABLE II Parameter Setup II

N	$\lambda$	$P_{on}$	T	k	f
10000	1	30%	10000	40	0.9

Figures 9 and 10 show the results of the number of overhead messages per query and stale query ratio versus different update frequencies. When the average update period ( $T_{up} = 10^5$ ) is much larger than the peer online and offline cycle, the overhead messages of the pull process are the major source. Due to the similar pull process, the number of overhead messages per query for two algorithms is close in this case. As the update period decreases, the number of overhead messages from the push process increases and dominates the number of overhead messages from the pull process. This leads to a better performance

of UPTReC than that of Rumor. When the update period is much shorter than the peer online and offline cycle, the number of overhead messages per query in UPTReC is more than 70% lower than that of the Rumor. The stale query ratio in UPTReC is less than 0.1% in all range of update periods. In Rumor, when the update frequency is high, the stale query ratio is about 2% for R = 80, and it is reduced to less than 1% when R = 100. The effect of R is similar to f. A larger R or f gives a lower stale query ratio but costs more overhead messages. The results show that the UPTReC can save up to 70% overhead messages while providing better probabilistically consistency guarantee for highly update files compared to the Rumor algorithm.



Fig. 9. The number of overhead messages per query versus average update period



Fig. 10. Stale query ratio (%) versus average update period

Through these comparisons, we know that the UPTReC can significantly reduce overhead messages to propagate an update with a smaller stale hit ratio comparing with Rumor algorithm.

#### VI. CONCLUSIONS AND FUTURE WORK

In this paper, we propose a novel algorithm, UPTReC, to propagate update through replica chain for decentralized and unstructured P2P systems. UPTReC provides probabilistically guaranteed file consistency. In UPTReC, each file has a logical replica chain composed of all RPs. Each RP has a partial knowledge of the chain. When an RP updates the file, it pushes the update to all possible online RPs through the replica chain. When an offline RP gets reconnected, the file status is synchronized by pulling an online RP.

An analytical model of the proposed algorithm is derived. The performance results of UPTReC compared to that of the Rumor algorithm shows that the UPTReC reduces up to 70% overhead messages to propagate updates with a smaller query ratio for highly updating files.

If each RP keeps a small number of probe peers, the update propagation may be stopped at some node. If a reconnected peer pulls an online peer from each direction and if an update is found, the peer can push the update to another direction of the chain. This process will reduce the stale query ratio by keeping small number of probe peers. The mechanism and the replica chain maintenance costs will be considered in our future work. As the growing of application in P2P systems, strong cache consistency is a further requirement, and this will also be considered in our future work.

#### REFERENCES

- [1] Open Source Community, Gnutella. In http://gnutella.wego.com, 2001
- [2] K. Aberer, P-Grid: A Self-organizing Access Structure for P2P information Systems. In *Proceedings of the Sixth International Conference on Cooperative Information Systems*, Trento, Italy, 2001.
- [3] K. Aberer, Z. Despotovic, Managing Trust in a P2P Information System. In Proceedings of the 10th International Conference on Information and Knowledge management, pp310-317, ACM press 2001.
- [4] R. Bhagwan, S. Savage and G. M. Voelker, Understanding Availability. In Proceedings of the 2nd International Workshop on Peer-to-peer systems, 2003.
- [5] Y. Chawathe, S. Ratnasamy, L. Breslau, N. lanham and S. Shenker. Making Gnutella-like Peer-to-Peer Systems Scalable. In *Proceedings of ACM SIGCOMM'03*, 2003.
- [6] E. Cohen and S. Shenker. Replication Strategies in Unstructured Peer-to-Peer Networks. In *Proceedings of the ACM SIGCOMM'02 Conference*, 2002.
- [7] E. Cohen, A. Fiat and H. Kaplan, Associative Search in Peer-to-Peer Networks: Harnessing and Latent Semantics. In *Proceedings of IEEE INFO-COM*'03, 2003.
- [8] A. Datta, M. Hauswirth and K. Aberer, Updates in Highly Unreliable, Replicated Peer-to-Peer Systems, In *Proceedings of IEEE ICDCS'03*, pp76-88, Rhode Island, May, 2003.
- [9] B. Gedik and L. Liu, PeerCQ: A Decentralized and Self-Configuration P2P Information Monitoring System, In *Proceedings of IEEE ICDCS'03*, pp490-499, Rhode Island, May, 2003.
- [10] S. Iyer, A. Rowstron and P. Druschel. Squirrel: A Decentralized Peer-topeer Web Cache. In Proceedings of the 21th ACM Symposium on Principles of Distributed Computing (PODC), 2002.
- [11] J. Lan, X. Liu, P. Shenoy and K. Ramaritham. Consistency Maintenance in Peer-to-Peer File Sharing Networks. In *Proceedings of Third IEEE* Workshop on Internet Applications, June, 2002
- [12] Q. Lv, P. Cao, E. Cohen, K. Li, and S. Shenker. Search and Replication in Unstructured Peer-to-Peer networks. In *Proceedings of the 16th Annual* ACM International Conference on Supercomputing, 2002
- [13] T. W. Page, R. G. Guly, J. S. Heidemann, D. Reiher, A. Goel, G. H. Kuenning, and G. J. Popek. Perspectives on Optimistically Replicated Peer-to-Peer Filing. *Software-Practice and Experience*, pp155-180, 28(2), 1998.
- [14] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content-Addressable Network. In *Proceedings of ACM SIGCOMM*, 2001
- [15] M. Roussopoulos and M. Baker. CUP: Controlled Update Propagation in Peer-to-Peer Networks. In *Proceedings of the 2003 Annual USENIX Technical Conference*, June 2003.
- [16] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan. Chord: A Scalable Peer-to-Peer Lookup Service for Internet Applications. In *Proceedings of ACM SIGCOMM*, 2001.