



Web search results caching service for structured P2P networks



Erika Rosas^b, Nicolas Hidalgo^{b,*}, Mauricio Marin^{a,b}, Veronica Gil-Costa^{a,c}

^a Yahoo! Labs Santiago, Chile

^b DIINF, University of Santiago, Chile

^c CONICET, National University of San Luis, Argentina

HIGHLIGHTS

- The paper proposes building a Web result cache on settop boxes available at homes.
- Settop boxes are organized as a P2P network with efficient algorithms for query routing.
- The P2P network contains algorithms to deal with peaks in user query traffic.
- The P2P network contains algorithms to increase cache hits in the user community.
- The proposed Web result cache reduces communication traffic outside the ISP network.

ARTICLE INFO

Article history:

Received 16 January 2013

Received in revised form

9 May 2013

Accepted 17 June 2013

Available online 8 July 2013

Keywords:

Web search engines

Caching services

Load balancing

P2P networks

ABSTRACT

This paper proposes a two-level P2P caching strategy for Web search queries. The design is suitable for a fully distributed service platform based on managed peer boxes (set-top-box or DSL/cable modem) located at the edge of the network, where both boxes and access bandwidth to those boxes are controlled and managed by an ISP provider. Our solution significantly reduces user query traffic going outside of the ISP provider to get query results from the respective Web search engine. Web users are usually very reactive to worldwide events which cause highly dynamic query traffic patterns leading to load imbalance across peers. Our solution contains a strategy to quickly ease imbalance on peers and spread communication flow among participating peers. Each peer maintains a local result cache used to keep the answers for queries originated in the peer itself and queries for which the peer is responsible for by contacting the Web search engine on-demand. When query traffic is predominantly routed to a few responsible peers our strategy replicates the role of “being responsible for” to neighboring peers so that they can absorb query traffic. This is a fairly slow and adaptive process that we call mid-term load balancing. To achieve a short-term fair distribution of queries we introduce a location cache in each peer which keeps pointers to peers that have already requested the same queries in the recent past. This lets these peers share their query answers with newly requesting peers. This process is fast as these popular queries are usually cached in the first DHT hop of a requesting peer which quickly tends to redistribute load among more and more peers.

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

Web search engines are systems devised to cope with a highly dynamic and demanding query rates, potentially of the order of many hundred thousand queries per second. To support redundancy and fault tolerance, large search engines operate on multiple, geographically distributed data centers. Intensity of query traffic is featured by the unpredictable behavior of users who are usually very reactive to worldwide events [1]. Web search queries typically follow a Zipf-like distribution of data [2,3] which produce load balancing problems over the query result cache servers.

* Corresponding author. Tel.: +56 29780671.

E-mail address: nicolas.hidalgo@usach.cl (N. Hidalgo).

This paper proposes implementing a result cache service for Web queries using the concept of a nano datacenter infrastructure implemented over a P2P network composed by small peers. The nano datacenter architecture uses ISP-controlled home gateways, like the ubiquitous set-top-boxes or STBs [4], to provide computing and storage services, and adopts a managed peer-to-peer model to form a distributed infrastructure [5].

We organize the STBs following a Distributed Hash Table (DHT) model which generally provides better routing/searching performance with much smaller overheads when compared against unstructured P2P overlays [6]. DHTs enable building scalable, decentralized and self-organizing systems. These systems balance storage, transparently tolerate peer failures, and provide efficient routing of queries. However, they suffer from load balancing problems under the presence of highly imbalanced query distributions. The emergence of hot-spots (popular queries) and flash crowds

can lead to saturation of peers degrading performance of a result caching service.

Caching schemes for DHTs and P2P networks have been proposed for applications such as Web caching and streaming, and have been the subject of research. In our context, the main difference with systems studied in previous works is that a cache service for Web search engine queries must be devised to be efficient under the following requirements: (1) query traffic can be very dynamic both in intensity and trending topics, and it can be subjected to sudden peaks occurring at unpredictable time instants; (2) cached queries are not expected to remain valid for a long time as current search engines supporting real time search may include new documents within a few minutes from publication on the Web; and (3) it is necessary to transfer the whole query answer to the requesting user as opposed to pointers to answer objects for later retrieval.

To address these requirements, a portion of each peer main memory is used as a result cache to hold query answers retrieved from one or more Web search engines. These queries are assumed to be distributed uniformly at random on N peers by means of a hash function on the query terms. They correspond to queries submitted by users in the recent past. Thus, any given peer is responsible for a fraction of these queries and (on request) it must contact the appropriate Web search engine to get the current answers for them. Users submit queries at their local peers, and each peer must respond with a query answer to its respective user.

The contribution of the proposed P2P result caching are mainly the following:

- A strategy that quickly reacts upon a sudden and drastic increase in query traffic where possibly query contents are highly skewed to a relative small number of different terms.
- A strategy to replicate the role of overloaded peers which provides load fairness among the peer replicas using a small amount of communication.

We propose using the two strategies in combination to efficiently cope with user query traffic on a P2P result caching system for Web search queries.

We evaluate our proposal by using a query log from a commercial Web search engine to compare it against state of the art baseline strategies. The results show that our approach (1) improves load balancing among participant peers, (2) increases cache hits, (3) reduces the amount of traffic generated within the peers community and towards the Web search engines, and (4) significantly reduces the communication volume associated with the replication of highly popular queries. All this leads to a more efficient performance than the alternative approaches, which also has a positive impact on reduction of overall power consumption.

An early version of this paper was presented in [7]. In this paper we extend [7] to significantly improve the performance of the strategy devised to react upon sudden query traffic peaks and present a more comprehensive experimental performance study.

The rest of this paper is organized as follows. Section 2 surveys related work and presents some background concepts about our solution. Section 3 presents the architecture, model, and the proposed solution. Section 4 presents a performance evaluation study conducted through process oriented discrete simulation. Finally, Section 5 presents concluding remarks.

2. Background and related work

In the following we briefly explain DHTs in order to make our paper self-contained. Then, we present related work on caching and replication in the context of DHTs.

2.1. Distributed hash tables

A Distributed Hash Table (DHT) is a self-organized structured substrate built over P2P overlay networks which provides data availability and persistence through replication. Every peer in such systems is associated with a unique identifier which defines the peer's position in the structure and the range of keys it is responsible for. Data is identified by a key in the same identifier space and can be located within a logarithmic number of routing hops.

DHT overlays maintain a strong topology, for example a ring in the case of Chord [8], Pastry [9] or Tapestry [10]. We design our solution in the context of Pastry [9] but it can be applied to other DHT realizations.

Every peer in Pastry is assigned a unique nodeID in a space of 128-bit identifiers generated using a cryptographic hash SHA-1. The neighbors of a peer in Pastry are stored in a *leafset* that contains the L numerically closest peers, $L/2$ clockwise and $L/2$ counter-clockwise. This set can handle replicas to improve fault tolerance in an environment where peers join and leave the network without warning. A “keep alive” message is used periodically to detect failed peers.

The Pastry routing algorithm is prefix based and it routes a message to the numerically closest peer of a given key k . In this paper we call this peer the *responsible* peer of k . The Pastry routing table stores on the n th row the IP address of peers whose nodeIDs share the first n digits with that of the present peer. The algorithm forwards the messages to a peer from its routing table that shares at least one more digit with the key k than the current peer. If no such a peer can be found and the current peer does not know any other peer numerically closer to k , then the current peer is the responsible of k and the routing ends.

In the case of Pastry, the *leafset* and the *routing table* compose the *out-links* of a peer.

2.2. Caching and replication in DHTs

Caching and replication in DHTs are active areas of research. P2P cache schemes have been proposed mostly for Web pages (Squirrel [11], BuddyWeb [12], Backslash [13]), storage systems (Aren [14,15]) and video streaming applications [16].

Approaches for cache mainly focus on optimizing one metric, such as hit rate, latency, or communication.

Kangasharju et al. [17] propose a caching P2P community, as our work does. Their key assumption is that traffic is cheaper within the community than transferring content from external nodes. They propose an algorithm called Top-K MFR that stands for most frequently requested objects. Each node tracks the files for which it has received a request and retrieves from the outside only the most requested ones. Their model assumes that cached objects are large in space, like videos, and that the size of space for caching is relatively small, namely there is space to hold in the order of dozen objects. This is why selecting the objects stored in cache is their main focus. Our solution, on the other hand, is designed to handle many small-sized objects which occupy space in the order of a few KBs.

Despotovic et al. [18] aim at improving caching with respect to the total traffic originated within the DHT search. They propose to replicate an object i when a high demand for i is detected by sending a reference of i to the previous hops that forwarded the queries for i in the past. A threshold is defined in order to trigger the replication of references to i . Unlike our work, this solution [18] replicates references to objects. Our solution builds on the need to send the full object (query results) to the requesting peer.

Tigelaar et al. [15] explore search result caching in a P2P model for information retrieval. They have found that a small size

cache offers performance comparable to an unbounded size cache. Moreover, using simulation they have shown that in distributed scenarios, caching can greatly reduce query load. Their results also show that the Least Recently Used (LRU) policy is the best approach in the P2P setting.

PCache [19] is a proactive caching strategy based on popularity of objects. Replicas are stored in the overlay nodes that are closest to the peer responsible for an object. It periodically estimates the popularity value of a content by using an aggregation protocol among the peers, which can be costly in a large scale network.

Some of the authors of this paper have proposed solutions for result caching outside the P2P model environment [20–22]. Their work has focused on caching results for Web search engines deployed on clusters of processors. We have used one of the terms presented in [20] to name one of the structures proposed in our paper: Location Cache. The Location Cache in [20] reduces the amount of hardware involved in the solution of search engine queries by storing the nodes which provided the results for a query in the backend. With this structure they can also select the search nodes most likely to provide a good approximated answer to queries so that queries are prevented from hitting all search nodes (processors) when operating under near saturation query loads.

Replication has been done mainly in two ways: *multiple hash strategy* and *path replication*. The multiple hash strategy is described in [23–25]. However, since DHTs exhibit path convergence, multiple hash strategies are less adapted to efficiently handle highly popular objects than path replication.

Most DHT based strategies select a group of neighbor peers to replicate the data (e.g., the leafset in Pastry [9] contains neighboring peers). This technique uses uniform replication in a fixed number of peers to achieve high availability and fault tolerance. In the path replication techniques, the goal is to have other peers to answer queries instead of the responsible peer. These replicas could be even closer to the source than the responsible peer.

Stading et al. [13] use a local cache diffusion method which pushes the replicas of a document (object) to one hop closer to the source of the last requests. This flood creates a *bubble* which grows in relation to the intensity of the flood, until no peer on the perimeter of the bubble observes high request rates. However, details on how this decision is taken are not given. They also propose a directory diffusion, but this is not adapted to hot-spots.

Bianchi et al. [26] compute a popularity value, maintaining a counter for each object and a query counter. When a peer is overloaded, it replicates the most popular objects in the most frequent last peer along the path towards the destination peer.

The solutions in [13,26] suffer from the *bubble* problem where a responsible peer stops answering queries for a popular key.

Yamamoto et al. [27] propose a path adaptive replication method which determines the probability of replication according to a predetermined replication ratio and storage capacity. Although it considers the storage availability of peers, it ignores their current load.

Ramasubramanian et al. [28] replicate objects based on the Zipf parameter α in order to minimize the resource consumption using local measurements and limited aggregation. However, flash crowds can deviate the requests from the Zipf distribution. A popular object is replicated by levels in all the peers that share one digit less with the key than the responsible peer. This work reduces latency of popular objects, but it creates a high number of replicas in the system which depends on the size of the network.

Silvestre et al. [14] propose a replication scheme called Aren, which takes into account content popularity and QoS metrics in content distribution systems for edge networks. Their strategy uses a hierarchical network for a cloud environment and relies on the use of a coordinator to optimize scheduling and replication. The use

of edge networks to build the system is similar to what is presented in this work. However, our P2P model is fully distributed.

In order to cope with load balancing issues in structured P2P networks, not only replication has been considered. Most of the solutions focus on providing a uniform distribution for the range of keys assigned to peers. Works like [29] for example, balance the namespace to distribute the workload evenly. The virtual server strategy is also widely used. In this case, a single physical node may have multiple virtual nodes in the logical structure. This model enables resource-aware load balancing as peers may vary the number of virtual nodes depending on their capacity. A drawback of this type of solution is that a physical node has to handle multiple routing tasks for maintenance and forwarding [30–32].

Solutions that modify a node's identifier to redistribute load do not consider the fact that query distribution can be highly skewed. If the identifier of a node changes, the assigned partition is moved between adjacent nodes. The same occurs if a node leaves and joins at a different place in the ring, like in [33,34]. The responsible peer changes but not the load hot-spot. Low loaded peers can exchange place with highly loaded peers [35,36], but the responsibility lies in a single peer. Moreover, security may be affected when allowing arbitrary node identifier modification as malicious peers could join the ring by adopting IDs to serve portions of queries. Ledlie et al. [37] propose using a strategy based on k -choices of safe IDs to solve this problem. Each trusted peer joining the network is allowed to choose k verifiable IDs which the peer can use to serve different seemingly random sectors of the DHT ring. Karger and Ruhl [38] propose that each node chooses $\log n$ places in the ring and initially takes the responsibility for only one of them. This scheme has been used together with the virtual server strategy.

In our work, as opposed to considering the distribution of content in the network, we focus on object replication since in our case we need to cope with hot spots and skewed distribution of queries.

3. Proposed P2P caching service

3.1. Architecture

We propose to build a Caching Service (CS) to reduce query traffic towards Web search engines by using a P2P overlay composed of peers within a single geographic region and operating as a *nano datacenter* infrastructure. The idea is to create a fully distributed service platform based on managed boxes (set-top-box or DSL/cable modem) located at the edge of the network, where both boxes and access bandwidth to those boxes are controlled and managed by a given entity (e.g., by Telco, virtual operator or service provider) [39]. Fig. 1 presents the proposed architecture. The set-top-boxes act like peers on the overlay. By utilizing virtualization technology the application running on a box can be completely transparent to the end user. The advantages of this approach rely on their low-energy consumption, self-scalability, and high availability. Additionally, a nano datacenter does not suffer from free-riding, peer dynamics, and lack of awareness of underlying network conditions [5]. Importantly, they provide access to a very large amount of resources. As an example, in USA there are more than 160 million set-top-boxes (STBs) which can potentially act as small servers [4].

For a P2P CS to be of practical interest the issue of low latency to get query answers is critical. On the other hand, structured P2P networks suffer from high latency because during search it is necessary to visit several peers to find the peer that contains the target object. In order to ease this problem, like in [11], we consider that peers remain within a geographic region.

In our architecture, user web browsers issue query requests to the P2P CS overlay. Peers (set-top boxes) route queries among

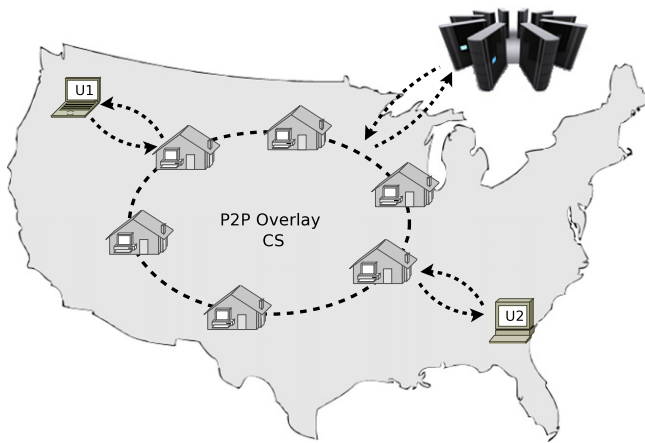


Fig. 1. Architecture.

peers in order to find cached information, namely answers of queries. Like in a standard DHT, terms of queries are hashed by using a cryptographic hash, such as SHA-1, and routed through the overlay to the peer which shares the closest ID with the hashed query. If the P2P CS finds a valid cached query answer it responds directly to the requesting user. Otherwise, the query is sent to the Web search engine to get the answer, cache it on the P2P CS and respond to the user. The main goal is to alleviate traffic towards search engines and reduce the costs for ISPs by reducing traffic outside their networks.

3.2. TLC+: Two-level caching

We assume a classical DHT overlay composed by N physical nodes (or peers) and K object-keys (query terms) mapped onto a ring. Objects (query terms plus the respective query answer) stored in a DHT such as Pastry or Chord, have a responsible peer in the network that is the peer with the closest ID to the key of the object. Thus, any given peer is responsible for a fraction of these queries and (on request) it must contact the appropriate Web search engine to get the current answers for them. Below we use query answer and query result to mean the same, namely a valid html page containing the results of executing the respective query at the Web search engine which can be directly sent to the user.

Participant peers have the ability to *receive*, *forward*, and *answer* a given query. In our model we assume that answering a query (i.e., sending a query answer to a requesting peer) consumes the highest amount of resources considering that query answers are expected to be of the order of tens of KBs. On the other hand, *receive* and *forward* operations are in the order of hundreds of bytes. Requesting an object to the Web search engine is assumed to be costly to the ISP since the communication traffic leaves its network.

Peers of the network create a collaborative cache which acts as a distributed CS. The set-top box provider defines a fixed storage space for caching query answers from the Web search engine for all peers.

The proposed P2P caching strategy is composed of two levels. The first one is devised to quickly react upon a sudden and drastic increase in query traffic where possibly query contents are highly skewed to a relatively small number of different terms. In this case, each peer p_i that submits a query q and receives the answer can potentially share it with other peers p_k that are requesting q during the same period of time. This quickly eases communication traffic on the peer r that is responsible for contacting the search engine and sending back the answer of query q to the requesting peer p_i . Namely, if peer p_1 requests the answer for a query q to its responsible peer r , and soon after the peers p_2 and p_3 request the answer for the same query q (in that order), then the proposed

strategy causes a query answer transfer from r to p_1 , then a transfer from p_1 to p_2 , and finally a transfer from p_2 to p_3 . Notice that the number of sequential transfers cannot be excessively large and parallelism can be exploited for a new peer p_4 as peer p_1 can transfer the query answer to p_4 in parallel with the transfer from p_2 to p_3 . Below we describe the sequential and parallel solutions separately. The former is suitable for moderate query traffic and uses minimum memory space whereas the later is suitable for high query traffic but it requires more memory space. We call this part short-term load balancing.

The second level, called mid-term load balancing, is in charge of replicating the role played by a peer r responsible of a query q in neighboring peers when traffic intensity exceeds a threshold value. In this case, peers that are neighbors to peer r in the sense of the DHT protocol, are also allowed to contact the Web search engine to get answers for the query q and pass them back to the requesting peers. The decision on what neighbors to replicate a query q is made by considering the intensity of query traffic arriving from each neighbor to peer r . Neighbors routing more instances of query q to peer r are more likely to get a replica of q , that is they are more likely to become a responsible peer for q .

3.3. Algorithmic details

Each peer has a standard result cache (RCache) and a special purpose cache called location cache (LCache). The RCache is used to store answers to queries and their respective query terms in a normalized format. This cache stores answers of queries for which the peer is responsible for (we call this subset as *global cache*), and answers of queries delivered to users that submitted queries in the peer (we call this subset as *local cache*). Each cache entry has a state flag indicating whether the respective query answer is being retrieved from a Web search engine or from some other peer in the network. Also, in each entry there is a timeout flag indicating the time instant from which the respective query answer is no longer valid. The timeout is determined from TTL values assigned to queries retrieved from the search engine (TTL stands for “time to live”). The RCache entries are administered with the LRU replacement policy.

The LCache stores data used to support short-term scheduling. Each entry stores terms of a query, an IP address of a peer that contains (or is to contain soon) the answer for the query, and an expiration counter used to limit the number of times the entry is used to route incoming queries with the same terms. These cache entries are also administered with the LRU replacement policy though their overall size is much smaller than in the RCache as they do not store query answers.

In the following we describe the combined operation of RCache and LCache in each peer through an example.

Assume that a peer d is responsible for a query q and that a peer a_1 requests the answer for q . Peer a_1 does not find an entry for q in either its RCache or its LCache so it must get the query answer from peer d . The DHT protocol establishes that to go from a_1 to d it is necessary to pass through peers b and c . In consequence, peer a_1 stores q in its RCache with flag state indicating “answer being retrieved” and it sends q to peer b . Since b is the first DHT hop from a_1 , an entry (q, a_1, m_b) is stored in the LCache of peer b where $m_b = m$ is the initial value for the entry expiration counter. Then peer b sends q to peer c , and as peer c does not find a valid entry for q in either its RCache or its LCache, it sends q to peer d . At this point, peer d searches for a valid entry for q in its RCache which, upon a cache miss, it triggers a request for q to the Web search engine. In the meantime, another peer a_2 generates a request for the same query q by sending q to peer b since b is the first DHT hop from a_2 as well. This causes a cache hit in the LCache of peer b which changes the entry associated with q from (q, a_1, m_b) to $(q, a_2, m_b - 1)$. Then

peer b sends the pair (q, c) to peer a_1 where the peer c address is included for cases in which peer a_1 contains an invalid entry for q in its RCache and q is not in its LCache either, so the query is sent from peer a_1 directly to peer c to reach peer d if necessary. In either case, an entry (q, a_2, m_a) is stored in the LCache of peer a_1 where $m_a = m_a - 1$ or $m_a = m$ depending on whether there is an entry for q already stored in the LCache or not. On the other hand, once peer d contains a valid answer for q in its RCache, it initiates a data transfer from peer d to peer a_1 to make it possible for peer a_1 to contain a valid answer for q stored in its RCache and respond to the requesting user. In turn, this triggers a data transfer from peer a_1 to peer a_2 so that a_2 can respond to its user and store the answer for q in its RCache.

In this way, the aim of the combined operation of the RCache and LCache objects in each peer is to quickly prevent the potential saturation of peer d when an extremely popular query q suddenly arises. Basically the main idea is to take advantage of the fact that it is necessary to transfer query answers to requesting peers which are re-used as replicas by other requesting peers. When the query traffic on peer d exceeds a threshold value, it shares the responsibility for q on one or more neighboring peers, such as peer c , to reduce the load on d . This peer c is then allowed to contact the search engine to let it hold a valid replica of the answer for query q . Overall, the LRU replacement policy takes care of cache entries that are no longer useful in the RCache and LCache objects.

Algorithms 1 and 2 describe the steps followed by the proposed two-level scheme upon the occurrence of events triggered by the reception of query messages in each peer.

Algorithm 1: Event handler in each peer

```

1 Event: A new query  $q$  arrives from a user connected to
2   this peer
3 begin
4   if RCache.find( query=  $q$  ) then
5     if RCache.isValid( query=  $q$  ) then
6       return RCache.answerQuery( query=  $q$  )
7   else if IsResponsible(  $q$  ) then
8      $r =$  getAnswerFromSearchEngine( );
9     RCache.insert( query=  $q$ , flag="retrieved", answer=  $r$  )
10    return  $r$ 
11  if LCache.find( query=  $q$  ) then
12     $p =$  LCache.getPeerIP( query=  $q$  )
13  else
14     $p =$  DHTgetNextPeerIP( thisPeer( ) )
15  RCache.insert( query=  $q$ , flag="retrieving" )
16  send( peer=  $p$ , query=  $q$ , source= thisPeer( ) )

17 Event: An answer for query  $q$  arrives from another peer
18 begin
19    $r =$  getAnswer(  $q$  )
20   RCache.update( query=  $q$ , flag="retrieved",
21     answer=  $r$  )
22   foreach query  $e$  waiting for  $r$  do
23     send( peer= sourcePeer(  $e$  ), query=  $e$ , answer=  $r$  )
24   return  $r$ 

```

3.4. Replication

Popular objects (or hot objects) may produce an overload for the responsible peer which may not be able to handle the amount of requests of a hot object. Our solution attempts to solve this issue by dynamically increasing the number of peers responsible for an object. The key idea is to spread the load among the replicas taking into account the current load and capacities of the peers. Unlike

Algorithm 2: Event handler in each peer (cont.)

```

1 Event: A query  $q$  arrives from another peer
2 begin
3   if RCache.find( query=  $q$  ) then
4     if RCache.isValid( query=  $q$  ) then
5        $r =$  RCache.answerQuery( query=  $q$  )
6       send( peer= sourcePeer(  $q$  ), query=  $q$ , answer=  $r$  )
7     else if containNextPeer(  $q$  ) then
8       send( peer= nextPeer(  $q$  ), query=  $q$  )
9   else if IsResponsible(  $q$  ) then
10     $r =$  getAnswerFromSearchEngine( );
11    RCache.insert( query=  $q$ , flag="retrieved", answer=  $r$  )
12    send( peer= sourcePeer(  $q$  ), query=  $q$ , answer=  $r$  )
13  else if LCache.find( query=  $q$  ) then
14     $p =$  LCache.getPeerIP( query=  $q$  )
15    LCache.update( query=  $q$ , peerIP= sourcePeer(  $q$  ) )
16     $c =$  DHTgetNextPeerIP( thisPeer( ) )
17    send( peer=  $p$ , query=  $q$ , nextPeer=  $c$  )
18  else
19    if firstHopDHT( sourcePeer( ) ) then
20      LCache.insert( query=  $q$ , peerIP= sourcePeer( ) )
21    if containNextPeer(  $q$  ) then
22      send( peer= extraPeer(  $q$  ), query=  $q$  )
23    else
24       $p =$  DHTgetNextPeerIP( thisPeer( ) )
25      send( peer=  $p$ , query=  $q$  )

```

previous works, we do not overload the saturated peer by sending the cached objects but sending a small message to authorize the peer to request the object directly from the Web search engine.

We define C_i the maximum capacity of a peer i , as the requests per unit of time it is able to answer, for example, 100 requests per second. When the number of queries received for a peer i exceeds C_i we say the peer is saturated. Each peer keeps track of the number of times an object k was accessed during a time interval. We call this value f_k , the frequency of k . The load of a peer is the total number of times it has sent a cached object to another peer, which also corresponds to the summation of the frequencies of all the objects it had in its cache in that interval.

In order to take replication decision the maximum capacity C_i and its current load L_i are sent periodically by piggybacking this information in the routing and keep-alive messages. Additionally, each peer i stores the query frequency f_{ij} coming from its in-links j for each stored object. An in-link j is a peer j that contains the address of peer i in its routing table and thereby j can directly send messages to i .

Algorithm 3 presents our replication approach. We select a set of in-links to replicate considering that (1) the summation of incoming frequencies for q has to be higher than the load it wants to alleviate, and (2) the total available capacity of the involved peers is enough to answer the number of requests. Notice that in case all in-links do not have enough capacity available, they have to be replicated further into their own in-links.

3.5. Parallel LCache transfers

For each query in the LCache we can reduce latency when a sudden peak in query traffic takes place. In this case, the number of requesting peers per unit time increases drastically. In such a situation letting sequential transfers of a query answer from one peer to the another can degrade performance since the last requesting peers should experience long latencies. We propose a solution for this case which consists of parallelizing query answer transfers.

Algorithm 3: Replication algorithm

```

input : Object  $q$ : most frequently accessed query in the peer
        RCache
input : Overload  $\beta$ 
1 sort ( $inLinks$ ): high to low  $f_{ij}$ 
2 while  $\beta > 0$  do
3    $P_i \leftarrow$  Next in-link
4   Replicate responsibility for  $q$  in  $P_i$ 
5   if  $f_{ij} < (C_i - L_i)$  then
6      $\beta = \beta - f_{ij}$ 
7   else
8      $\beta = \beta - (C_i - L_i)$ 

```

For a highly popular query in the LCache we can keep more than one peer ID, say D peer IDs. Thus, for a suddenly popular query q , the first requesting peer p_1 would have to get the query answer from the responsible peer and place its ID in the LCache entry associated with q . The second peer p_2 gets the answer from peer p_1 . After this transfer, both peer p_1 and p_2 are in condition to transfer the q answer to a third peer p_3 . Assuming that a transfer takes τ units of time, on average, peer p_3 should expect to receive the q answer no earlier than $2 \cdot \tau$ either from p_1 or p_2 . For a continuous stream of newly requesting peers, the number of peers available for transfers increases significantly. For instance, for $D = 2$ the number of available peers follows a Fibonacci sequence for the number of available peers able to deliver at $\tau, 2 \cdot \tau, \dots, n \cdot \tau$ units of time respectively.

On the other hand, keeping D unbounded does not increase performance after a threshold value since the peer arrival rate λ limits the total number of peers that can be served in parallel. In terms of a $G/G/\infty$ queuing model of this problem, the average number of peers being served at any time is given by $\lambda \cdot \tau$ and thereby, at steady state, it does not make sense to have $D > \lambda \cdot \tau$ in the LCache for q . This indicates the basis of our approach.

We let the value of D grow and shrink in accordance with the observed arrival rate λ . We keep a priority queue to retrieve the peer ID which is able to serve a transfer in the least time.

Algorithm 4 presents our strategy for selecting the next peer to let parallel transfer at the lowest cost in accordance with the current value of D .

Algorithm 4: GetPeer

```

1 LCache request: A new query  $q$  arrives from peer  $R$ 
2 begin
3    $t = \text{current\_time}()$ 
4    $\lambda = \text{update\_rate}(t)$ 
5    $n = \text{heap.size}()$ 
6   if  $n = 0$  then
7      $\text{heap.insert}(\text{new}(R), t + \tau)$ 
8     return NULL
9    $D = \lambda \cdot \tau$ 
10   $p = \text{heap.get}()$ 
11   $t_p = p.\text{timeout}()$ 
12  if  $t_p > t$  then
13     $t = t_p$ 
14  if  $n < D$  then
15     $\text{heap.insert}(\text{new}(R), t + \tau)$ 
16  if  $n + 1 < D$  then
17     $\text{heap.insert}(p)$ 
18  return  $p.\text{ID}$ 

```

4. Evaluation

We have used a process-driven simulation to conduct our evaluation using the simulation library libcppsim [40]. Libcppsim

Table 1

Default configuration.

Parameter	Value
Network size (num. peers)	1000
RCache size (num. entries)	100
Local cache size	50
LCache size (num. entries)	50

is a general purpose library written in C++ where each process is implemented by a co-routine that can be locked and unlocked during simulation. We have built a simulator that implements a transport layer, a P2P overlay and our caching proposal. Pastry [9] has been used as the overlay in our experimentation. We have also simulated the Web search engine process, which sends the answers to the responsible peers and other users that may exist outside the P2P community.

In our previous work [7] we used an event-driven simulation over Peersim [41]. However we were not able to have a dynamic query rate with this simulator. In the present work, we therefore use our own simulator to create a more realistic environment with different query rates to simulate flash crowd queries.

From the state of the art we have chosen two efficient strategies for DHTs to be compared against our approach. The first one, that we call *Leafset*, replicates objects of an overloaded peer in its closest neighbors as suggested in [9] to avoid saturation. The second approach, that we call *Bubble*, replicates objects of an overloaded peer in all of its in-links [18,13]. In order to make a fair comparison, we also included a Local cache to both approaches. These two approaches represent the main ideas from the state of the art.

The approach first presented in [7] is named TLC in figures below and the extended version of TLC proposed in this paper is called TLC+.

The default configuration of our simulation experiments is presented in Table 1. The simulation is divided in simulation time windows of 100 units of time. In our default configuration the query rate is 1000 queries per unit of time, and every 500 units of time we increase the query rate to simulate a flash crowd introducing 5000 queries for 5 units of time.

We have considered the freshness of the query answers by assuming that queries are assigned a time to live (TTL) value by the Web search engine (WSE) after which the cached object become stale. The TTL assigned by the WSE is independent of our cache solution. As in [42], we use the incremental approach to estimate the TTL value for objects cached by the responsible peers. Notice that the WSE does not reveal the TTL values assigned to query answers and thereby peers can only predict them.

Queries were extracted from a 16 million US–Canada user query log of the Yahoo! search engine. In the following we present our performance results.

4.1. Load balancing

In order to measure load balance among peers we use a metric based on Lorenz curves called the Gini coefficient, which is a metric commonly used in other fields like economics and ecology.

If all peers have the same load, the Lorenz curve is a straight diagonal line, called the line of equality or perfect load balancing. If there is any imbalance, then the Lorenz curve falls below the line of uniformity. The total amount of load imbalance can be summarized by the Gini coefficient G , which is defined as the relative mean difference, i.e. the mean of the difference between every possible pair of peers, divided by their mean load. For n peers the Gini coefficient is calculated by (1), where l_i correspond to the load of the peer i and μ is the mean load. G values range from 1 to 0. The value 0 is achieved when all peers have the same load. The value 1 is achieved when one peer receives the whole system load while

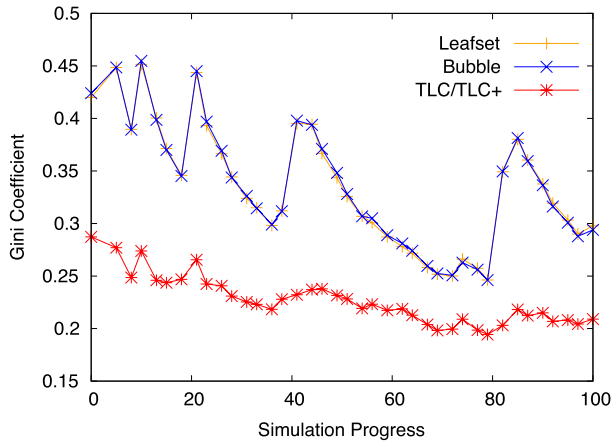


Fig. 2. Gini coefficient measuring load balance.

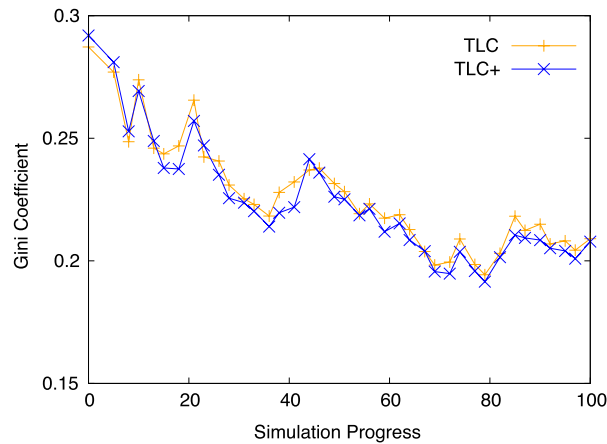


Fig. 3. TLC versus TLC+.

the remaining peers receive none. Therefore when G approaches 0, global load imbalance is small, and when G approaches 1 the imbalance is large.

$$G = \frac{\sum_{i=1}^n (2i - n - 1) \cdot l_i}{n^2 \cdot \mu}. \quad (1)$$

Notice that a Gini value equal to 0 is extremely rare in a P2P network. In [23] Pitoura et al. estimate that Gini values under 0.5 represent a fair load imbalance among the peers.

Fig. 2 presents a comparison of load balance conditions in the network by measuring the Gini coefficient for the TLC+, TLC, Leafset and Bubble strategies. The results show that our approach achieves better load balance than the other approaches. This is explained by our short-term scheduling. This strategy improves the hit rate of the local cache, since there is one node strategically chosen that knows its location. We reach an average improvement of 31% during the simulation as compared to the other techniques. Fig. 3 shows the difference between TLC+ and TLC. Both achieve very similar results. There is a 1% of improvement of TLC+ compared to TLC with respect to overall network load balance.

Flash crowds can be clearly seen in the Gini coefficient curves of the other approaches since the imbalance increases drastically. In our approach, on the other hand, flash crowds do not affect performance significantly.

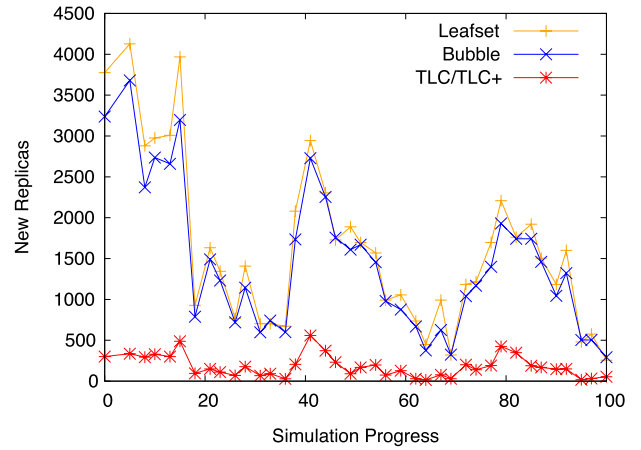


Fig. 4. Number of replicas.

4.2. Number of replicas

Another important metric is the number of replicas created in order to cope with peaks in load. The replication of an object in another peer implies the transfer of the whole object and when a peer is already overloaded; this can degrade the system performance even more. TLC/TLC+ only replicates the peer role of being the responsible for an object. This requires a low cost message but increases traffic outside the network.

Fig. 4 compares the number of replicas created during the evolution of the simulation. The number of replicas for TLC/TLC+ is significantly smaller than the number of replicas in the other approaches. This improvement is explained by our short-term scheduling strategy. The LCache greatly reduces the load on the responsible peers by redirecting requests to peers that have recently requested the queries and have already received the results. The peaks observed in the number of replicas generated by the other approaches are a consequence of flash crowd queries, and the new versions of expired-TTL objects that these strategies have to replicate in other peers.

The Leafset and Bubble replication schemes do not consider the load of other peers when responsible peers send objects to them. However, if they use the same scheme as the TLC/TLC+ strategy, that is, they perform replication by sharing the responsibility of requesting objects to the WSE among several peers, they could decrease the number of replicas as shown in Fig. 5. Nevertheless, in this case, the TLC/TLC+ strategy still produces less replication than the Leafset and Bubble strategies.

The extension of the LCache structure of TLC to support more than one peer ID (TLC+) has a positive impact on the number of replicas. The simulations assume a transfer latency between peers of 0.05 (low latency) and 0.5 (high latency) time units. Fig. 6 shows that the number of replicas decreases even more with the use of TLC+. With this solution more traffic is handled by the peer LCache structures deviating load to the objects located in the peer local caches.

4.3. Traffic generated

We analyze the traffic generated within the P2P network to respond to a query and the traffic generated towards the WSE by measuring the cache hits and misses respectively.

A cache hit takes place when an object (valid query answer) is found in the P2P network and thereby there is no need to fetch the object from the WSE. Fig. 7 presents the evolution of cache hits through simulation progress. TLC and TLC+ equally improve cache hits, thus, reducing in about 23% the amount of traffic generated

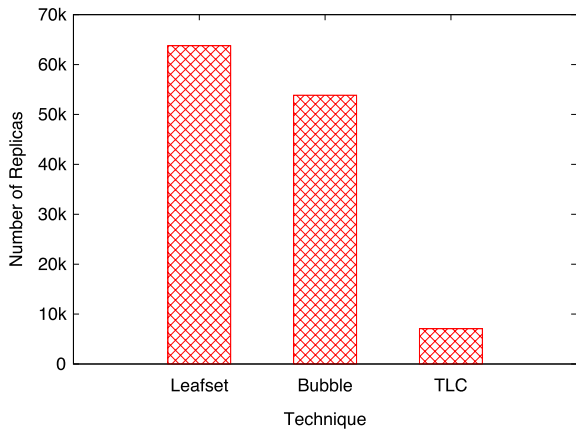


Fig. 5. Number of responsibility replicas.

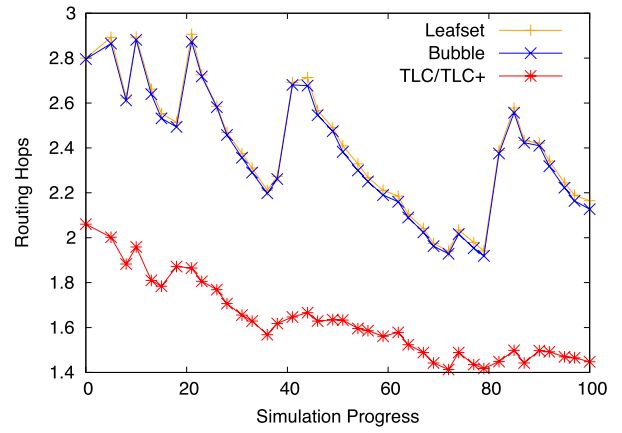


Fig. 8. Average number of hops.

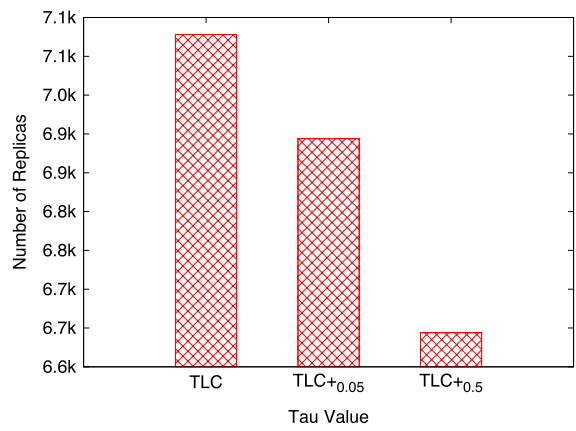


Fig. 6. Number of replicas $D > 1$.

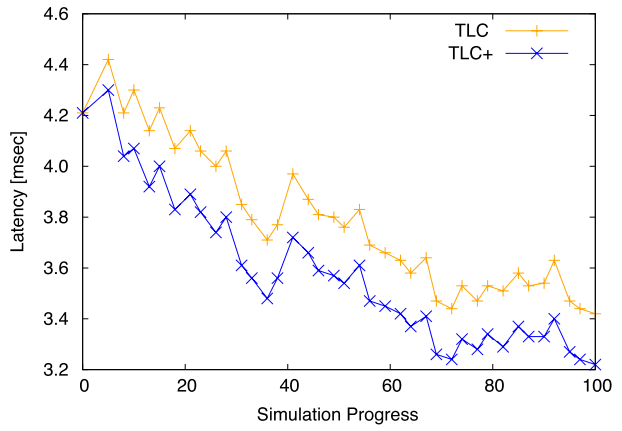


Fig. 9. Latency TLC versus TLC+.

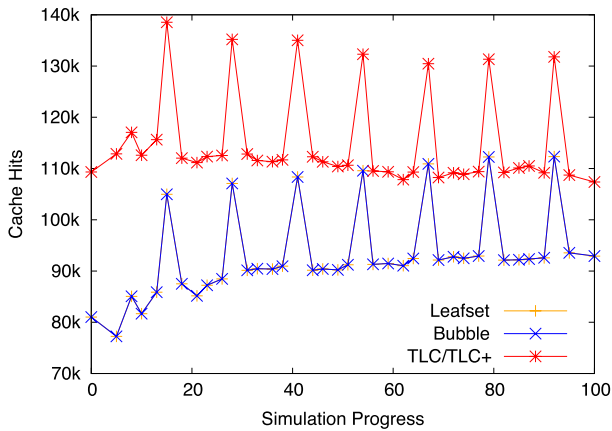


Fig. 7. Number of hits.

towards the WSE as compared against the alternative approaches. Leafset and Bubble also achieve the same number of hits. Peaks in the curves are a consequence of flash crowd queries we injected during the simulations.

4.4. Latency

We have also measured latency in terms of the number of hops required to find a cached object in the P2P network. Fig. 8 presents the results. TLC/TLC+ reduces the number of hops due to our short-term scheduling which quickly redirects popular queries to peers where it is likely to find the requested object in a local cache.

In the beginning of the simulation the number of hops starts at a higher point since cache structures are empty. During the simulation, the curves tend to decrease as the local caches start to contain valid answers for popular queries and therefore the routing process does not reach the responsible peers for all user searches.

The Leafset and Bubble approaches are greatly affected by TTL expiration since they cannot use invalid query answer copies stored in local caches and therefore searches have to reach the responsible peers. Our approach, on the other hand, copes smoothly with TTL expiration and maintains a stable number of hops during the whole search process.

TLC and TLC+ have same results in terms of routing hops. However, the TLC+ strategy decreases average response time even for a pessimistic case in which communication between peers is so fast that there is little margin to exploit the object transfer parallelism introduced by TLC+. Fig. 9 presents the query latency improvement introduced by TLC+ for the pessimistic case. TLC+ enables faster searches than TLC since more peers can benefit from the newly arriving copies of objects stored in the peer local caches at the same time.

4.5. Cache size

Fig. 10(a) and (b) show results when simulating different cache size configurations for the global and local cache entries present in the RCache, and the location cache entries present in the LCache. In each plot we show two measures: (1) the number of overall cache hits in the left y axis, and (2) the number of replicas in the right y axis. The x axes show the number of entries in each type of cache. In both plots, a hit means a query result that is found in some peer

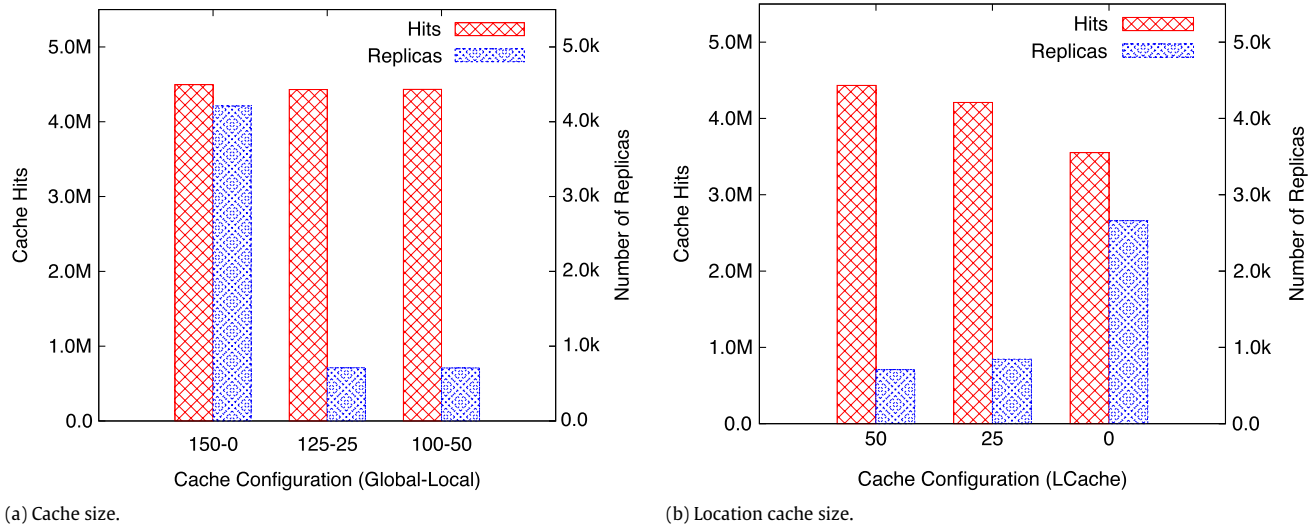


Fig. 10. Cache structures size.

RCache, and a replica means the assignment of a new DHT neighboring peer to share the responsibility of contacting the WSE for a given popular query.

The results in Fig. 10(a) show that the introduction of a small percentage of local cache entries in the RCache greatly reduces the number of replicas in the system. We did not find much difference between introducing 16% and 33% of local cache entries in the fixed space assigned to keep the RCache in each peer. In fact, the figure shows a slight decrease in hits when we increase the local cache entries whereas replication remains fairly constant. This is because the number of entries that are no longer available to the whole P2P community is increased when the global cache entries are reduced. Popular entries in the local cache, on the other hand, can become distributed in other peers in the system so they can be re-used by closer peers in the DHT. Thus the results indicate that it is sufficient to assign a small percentage (16%) of local cache entries in the space available for caching query results in each peer.

The results of Fig. 10(a) were obtained with all LCaches containing 50 entries. We have also analyzed the effects of different sizes for the LCaches. Fig. 10(b) presents how hit and replication performance degrade as we reduce the number of entries in the LCaches. Without the LCache entries, the popular query answers in the local caches are not found and thereby most of the searches arrive at the responsible peers. This overloads the responsible peers and consequently more replicas are generated to ease query load.

The results in Fig. 10(b) also show that even with a fairly small size for the LCaches (i.e., 50 entries), it is possible to decrease the number of replicas and increase the number of hits significantly. Note that each LCache entry is much smaller in space than a corresponding RCache entry. Apart from the query terms, the LCache entries mainly contain peer IDs whereas RCache entries contain full query answers. The space occupied by the LCache is less than 1% of the whole space reserved for caching in peers.

4.6. Power consumption

To evaluate our proposal against the alternative P2P caching strategies in terms of power consumption, we extend our simulators to measure hardware utilization and through this metric we estimate the energy required to process large sets of Web queries in the network. We set the simulator energy consumption curves in accordance with the utilization-to-energy curves presented in [43] for different commodity hardware devices. The effects of varying

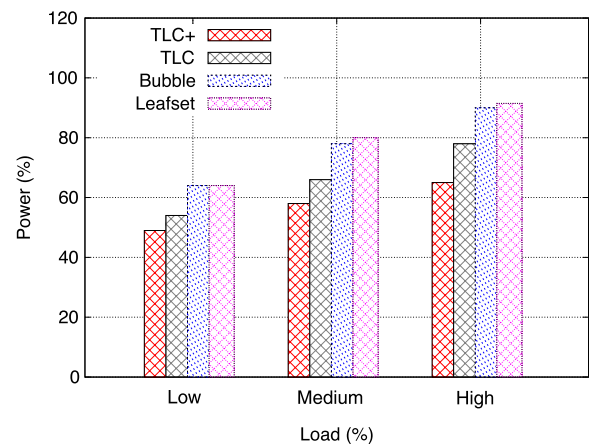


Fig. 11. Power consumption results.

query traffic intensity and flash crowd queries are reflected in the level of utilization of devices which has a non-linear corresponding effect in energy consumption. Typically at 0% utilization devices still consume more than 50% of peak power. To set power values close to the set-top boxes and network setting we use the baseline values presented in [5].

Fig. 11 shows the impact on reduction of power consumption of the proposed TLC and TLC+ strategies under different workloads as compared against the alternative strategies. The results show that the better efficiency of the proposed TLC+ strategy, in terms of communication and overall cache hits, can lead to a significant reduction of total power consumption with respect to the other cache strategies.

Note that the work in [5] makes the case for nano data centers by studying video-on-demand applications. They conclude that, even under pessimistic scenarios, overall energy consumption may be reduced between 20% and 30% when compared against centralized data centers. They argue that the additional energy cost for home users is not significant. The small extra cost in electricity bills can be amortized by the ISP with service discounts and incentives for home users.

As we explain in the following, our application requires much less computing and communication than the video-on-demand application and thereby it consumes less energy. As to communicating single objects inside the network, the size of a query answer is of the order of a few KBs whereas video sizes are of

the order of MBs and GBs. Current trends in Internet usage indicate that Web and E-mail applications represent less than 20% of the total content traffic [5]. Thus, at steady state, user query traffic should consume a small fraction of network resources. On the other hand, computing is also expected to be very small because (1) we perform a few low-cost operations per query on caches and (2) it is possible to achieve a significant improvement in performance with very few cache entries per peer. The latter saves energy as there is no need to use hard-disks in set-top boxes as it may be the case for the video-on-demand application. Therefore, the effect on home electricity bills should be negligible.

5. Conclusions

We have proposed a result cache strategy for structured P2P networks called Two-Level Result Caching (TLC+). It applies a short-term load balancing strategy based on the use of a LCache structure in each peer. The LCache is a tiny LRU cache that stores information about the peers that have recently requested queries in their first hop of the DHT routing path. The LCache increases parallelism in the network by making available to other peers the queries stored in the local result caches of peers that recently requested the same queries. In addition, a mid-term load balancing strategy is applied to replicate the role of being a responsible peer in its DHT neighboring peers. Responsible peers are in charge of contacting third-party Web search engines to get answers to subsets of user queries. They store these answers in their global result caches for future references. The replication strategy applied in the mid-term strategy takes into consideration the peer capacity and its current load in order to decide when to replicate responsibility for very popular queries.

The experimental results confirm that our solution improves load fairness among peers by 5% as compared to the best performing alternative strategy evaluated. Query answer replication is reduced up to 96% by our solution as compared to alternative strategies. Replica generation is a communication bandwidth consuming operation which can saturate peers and increase the cache miss rate degrading overall performance. Moreover, TLC+ improves in about 10% the number of cache hits as compared against the best alternative strategy, and about 18% as compared against the worst performer. This implies that traffic towards the Web search engine can be reduced by a similar percentage.

Our experimentation contemplated flash crowd queries, namely situations in which a sudden peak of a few popular queries are introduced in the simulation. This allows the evaluation of real-life situations where users become reactive to big news events. This is a type of experiment that has not been considered by related work and we have observed that it can degrade performance significantly. In addition, our experiments considered the execution of millions of queries submitted by actual users of a commercial Web search engine. The results show that TLC+ is more robust to query peaks than the alternative strategies as it can quickly spread the load generated by these queries among several peers. Overall, the proposed TLC+ improves load balance across peers, reduces object replications, increases cache hits, reduces the number of network hops, reduces individual query latency, requires a small number of cache entries per peer to be efficient, and demands less power consumption than alternative P2P caching strategies.

References

- [1] V.G. Costa, J. Lobos, A. Inostroza-Pisaj, M. Marín, Capacity planning for vertical search engines: an approach based on coloured petri nets, in: PETRI NETS 2012, in: Lecture Notes in Computer Science, vol. 7347, Springer, 2012, pp. 288–307.
- [2] L. Breslau, P. Cao, L. Fan, G. Phillips, S. Shenker, Web caching and zipf-like distributions: evidence and implications, in: IEEE INFOCOM'99, Vol. 1, 1999, pp. 126–134. <http://dx.doi.org/10.1109/INFCOM.1999.749260>.

- [3] K. Gummadi, R. Dunn, S. Saroiu, S. Gribble, H. Levy, J. Zahorjan, Measurement, modeling, and analysis of a peer-to-peer file-sharing workload, in: SOSP, 2003, pp. 314–329.
- [4] N.R.D. Council, Improving the efficiency of television set-top boxes. <http://www.nrdc.org/energy/files/settopboxes.pdf>.
- [5] V. Valancius, N. Laoutaris, L. Massoulié, C. Diot, P. Rodriguez, Greening the Internet with nano data centers, in: CoNEXT'09, ACM, 2009, pp. 37–48.
- [6] J.H. Ahn, U. Lee, H.J. Moon, Geoserv: a distributed urban sensing platform, in: CCGRID, 2011, pp. 164–173.
- [7] N.H. Erika Rosas, M. Marín, Two-level result caching for web search queries on structured p2p networks, in: Proceedings of the 2012 IEEE 18th International Conference on Parallel and Distributed Systems, ICPADS'12, IEEE Computer Society, Washington, DC, USA, 2012, pp. 221–228. <http://dx.doi.org/10.1109/ICPADS.2012.39>.
- [8] I. Stoica, R. Morris, D.R. Karger, M.F. Kaashoek, H. Balakrishnan, Chord: a scalable peer-to-peer lookup service for Internet applications, in: SIGCOMM, 2001, pp. 149–160.
- [9] A. Rowstron, P. Druschel, Pastry: scalable, decentralized object location, and routing for large-scale peer-to-peer systems, in: Middleware 2001, in: LNCS, vol. 2218, Springer-Verlag, 2001, pp. 329–350.
- [10] B.Y. Zhao, L. Huang, J. Stribling, S.C. Rhea, A.D. Joseph, J. Kubiatowicz, Tapestry: a resilient global-scale overlay for service deployment, IEEE Journal on Selected Areas in Communications 22 (1) (2004) 41–53.
- [11] S. Iyer, A.I.T. Rowstron, P. Druschel, Squirrel: a decentralized peer-to-peer web cache, in: PODC, 2002, pp. 213–222.
- [12] X. Wang, W.S. Ng, B.C. Ooi, K.-L. Tan, A. Zhou, Buddyweb: a p2p-based collaborative web caching system, in: NETWORKING 2002 Workshops, in: LNCS, vol. 2376, Springer, 2002, pp. 247–251.
- [13] T. Stading, P. Maniatis, M. Baker, Peer-to-peer caching schemes to address flash crowds, in: IPTPS'01, Springer-Verlag, 2002, pp. 203–213.
- [14] G. Silvestre, S. Monnet, R. Krishnaswamy, P. Sens, AREN: a popularity aware replication scheme for cloud storage, in: Proceedings of the 2012 IEEE 18th International Conference on Parallel and Distributed Systems, ICPADS'12, IEEE Computer Society, Washington, DC, USA, 2012, pp. 189–196.
- [15] A.S. Tigelaar, D. Hiemstra, D. Trieschnigg, Search result caching in peer-to-peer information retrieval networks, in: A. Hanbury, A. Rauber, A.P. de Vries (Eds.), IRFC, in: Lecture Notes in Computer Science, vol. 6653, Springer, 2011, pp. 134–148.
- [16] T. Fujimoto, R. Endo, K. Matsumoto, H. Shigeno, Video-popularity-based caching scheme for p2p video-on-demand streaming, in: AINA'11, IEEE Computer Society, 2011, pp. 748–755.
- [17] J. Kangasharju, K.W. Ross, D.A. Turner, Adaptive content management in structured p2p communities, in: InfoScale'06, ACM, 2006, p. 24.
- [18] Z. Despotovic, Q. Hofstätter, M. Michel, W. Kellerer, An operator approach to popularity-based caching in dhts, in: ICC, IEEE, 2010, pp. 1–6.
- [19] W. Rao, L. Chen, A.-C. Fu, G. Wang, Optimal resource placement in structured peer-to-peer networks, IEEE Transactions on Parallel and Distributed Systems 21 (7) (2010) 1011–1026. <http://dx.doi.org/10.1109/TPDS.2009.136>.
- [20] F. Ferrarotti, M. Marín, M. Mendoza, A last-resort semantic cache for web queries, in: J. Karlgren, J. Tarhio, H. Hyvärinen (Eds.), SPIRE, in: Lecture Notes in Computer Science, vol. 5721, Springer, 2009, pp. 310–321.
- [21] M. Marín, F. Ferrarotti, M. Mendoza, C. Gómez-Pantoja, V.G. Costa, Location cache for web queries, in: D. W.-L. Cheung, I.-Y. Song, W.W. Chu, X. Hu, J.J. Lin (Eds.), CIKM, ACM, 2009, pp. 1995–1998.
- [22] M. Marín, V.G. Costa, C. Gómez-Pantoja, New caching techniques for web search engines, in: S. Hariri, K. Keahey (Eds.), HPDC, ACM, 2010, pp. 215–226.
- [23] T. Pitoura, N. Ntarmos, P. Triantafyllou, Replication, load balancing and efficient range query processing in dhts, in: EDBT, 2006, pp. 131–148.
- [24] G. Swart, Spreading the load using consistent hashing: a preliminary report, in: SPDC'04, IEEE Computer Society, 2004, pp. 169–176.
- [25] D. Bauer, P. Hurley, M. Waldvogel, Replica placement and location using distributed hash tables, in: IEEE LCN 2007, IEEE Computer Society, 2007, pp. 315–324.
- [26] S. Bianchi, S. Serbu, P. Felber, P. Kropf, Adaptive load balancing for dht lookups, in: ICCN 2006, 2006, pp. 411–418.
- [27] H. Yamamoto, D. Maruta, Y. Oie, Replication methods for load balancing on distributed storages in p2p networks, in: SAINT'05, IEEE Computer Society, 2005, pp. 264–271.
- [28] V. Ramasubramanian, E.G. Sirer, Beehive: O(1) lookup performance for power-law query distributions in peer-to-peer overlays, in: NSDI, 2004, pp. 99–112.
- [29] X. Wang, Y. Zhang, X. Li, D. Loguinov, On zone-balancing of peer-to-peer networks: analysis of random node join, SIGMETRICS Performance Evaluation Review 32 (2004) 211–222.
- [30] B. Godfrey, I. Stoica, Heterogeneity and load balance in distributed hash tables, in: INFOCOM, 2005, pp. 596–606.
- [31] Y. Zhu, Y. Hu, Efficient, proximity-aware load balancing for dht-based p2p systems, IEEE Transactions Parallel and Distributed Systems 16 (2005) 349–361. <http://dx.doi.org/10.1109/TPDS.2005.46>.
- [32] H.-C. Hsiao, H. Liao, S.-T. Chen, K.-C. Huang, Load balance with imperfect information in structured peer-to-peer systems, IEEE Transactions on Parallel and Distributed Systems 22 (4) (2011) 634–649. <http://dx.doi.org/10.1109/TPDS.2010.105>.
- [33] M. Bienkowski, M. Korzeniowski, Friedhelm, dynamic load balancing in distributed hash tables, in: IPTPS, 2005, pp. 217–225. URL <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.59.1933>.
- [34] G. Giakkoupis, V. Hadzilacos, A scheme for load balancing in heterogeneous distributed hash tables, in: Proceedings of the Twenty-fourth Annual ACM Symposium on Principles of Distributed Computing, PODC'05, ACM, New York, NY, USA, 2005, pp. 302–311.

- [35] C. Chen, K.-C. Tsai, The server reassignment problem for load balancing in structured p2p systems, *IEEE Transactions on Parallel and Distributed Systems* 19 (2) (2008) 234–246. <http://dx.doi.org/10.1109/TPDS.2007.70735>.
- [36] H. Shen, C.-Z. Xu, Locality-aware and churn-resilient load-balancing algorithms in structured peer-to-peer networks, *IEEE Transactions on Parallel and Distributed Systems* 18 (6) (2007) 849–862. <http://dx.doi.org/10.1109/TPDS.2007.1040>.
- [37] J. Ledlie, M.I. Seltzer, Distributed, secure load balancing with skew, heterogeneity and churn, in: *INFOCOM*, IEEE, 2005, pp. 1419–1430.
- [38] D.R. Karger, M. Ruhl, Simple efficient load balancing algorithms for peer-to-peer systems, in: *Proceedings of the Sixteenth Annual ACM Symposium on Parallelism in Algorithms and Architectures*, SPAA'04, ACM, New York, NY, USA, 2004, pp. 36–43.
- [39] N. Laoutaris, P. Rodriguez, L. Massoulié, Echos: edge capacity hosting overlays of nano data centers, *SIGCOMM Computer Communication Review* 38 (1) (2008) 51–54. <http://dx.doi.org/10.1145/1341431.1341442>. URL <http://doi.acm.org/10.1145/1341431.1341442>.
- [40] M. Marzolla, Libcppsim: a Simula-like, portable process-oriented simulation library in C++, in: *ESM*, 2004, pp. 222–227.
- [41] M. Jelasity, A. Montresor, G.P. Jesi, S. Voulgaris, The Peersim simulator. <http://peersim.sf.net>.
- [42] S. Alici, I.S. Altingovde, R. Ozcan, B.B. Cambazoglu, O. Ulusoy, Adaptive time-to-live strategies for query result caching in web search engines, in: *Proceedings of the 34th European Conference on Advances in Information Retrieval*, ECIR'12, Springer-Verlag, Berlin, Heidelberg, 2012, pp. 401–412.
- [43] L.A. Barroso, U. Hölzle, The case for energy-proportional computing, *Computer* 40 (12) (2007) 33–37.



Nicolas Hidalgo obtained his Ph.D. in computer science from the Pierre and Marie Curie University, Paris, France. He is currently assistant professor at Universidad de Santiago de Chile. His research areas cover distributed systems, Peer-to-Peer resource discovery methods and complex queries over DHT-based networks. He has worked in localization mechanisms as well.



Mauricio Marin is a professor of computer engineering at University of Santiago, Chile. His Ph.D. is from the University of Oxford, UK (1999) with a thesis dissertation on parallel computing. Currently he is the head of Yahoo! Labs Santiago in Chile which is a center for applied research on scalable infrastructure for the Web and information retrieval. His research work is on optimizations for Web search engines.



Erika Rosas obtained her Ph.D. in computer science from the Pierre and Marie Curie University in Paris, France. She is currently assistant professor at Universidad de Santiago de Chile. Her research areas cover large scale networks, such as P2P and social network, in areas like trust and reputation systems.



Veronica Gil-Costa received her M.S. (2006) and Ph.D. (2009) in Computer Science, both from Universidad Nacional de San Luis (UNSL), Argentina. She is currently a professor at the University of San Luis, Assistant Researcher at the National Research Council (CONICET) of Argentina and a researcher for Yahoo! Labs Santiago (YLS). Her research interests are in the field of performance evaluation, similarity search, and distributed computing.