

Popularity-aware Prefetch in P2P Range Caching

Qiang Wang
University of Waterloo
q6wang@uwaterloo.ca

Khuzaima Daudjee
University of Waterloo
kdaudjee@uwaterloo.ca

M. Tamer Özsu
University of Waterloo
tozsu@uwaterloo.ca

Abstract

Unstructured peer-to-peer infrastructure has been widely employed to support large-scale distributed applications. Many of these applications, such as location-based services and multimedia content distribution, require the support of range selection queries. Under the widely-adopted query shipping protocols, the cost of query processing is affected by the number of result copies or replicas in the system. Since range queries can return results that include poorly-replicated data items, the cost of these queries is usually dominated by the retrieval cost of these data items. In this work, we propose a popularity-aware prefetch-based approach that can effectively facilitate the caching of poorly-replicated data items that are potentially requested in subsequent range queries, resulting in substantial cost savings. We prove that the performance of retrieving poorly-replicated data items is guaranteed to improve under an increasing query load. Extensive experiments show that the overall range query processing cost decreases significantly under various query load settings.

1 Introduction

Peer-to-Peer (P2P) infrastructure is being widely employed to support large-scale distributed applications. Besides well-established keyword-based file sharing and content distribution systems (e.g., Gnutella¹ and BitTorrent²), many applications require complicated query types such as range selection queries³ to be supported.

Range queries are used to retrieve all data that satisfy the specified *range constraints*. For instance, in on-demand P2P video systems (e.g., Joost⁴), video clip data within a certain time frame are buffered and shared by some peers, which can be modeled as range data. Peers can pose a range query over any time frame to search the video clips that they intend to play, and the video data satisfying the range constraints are returned to the query issuer. Similarly, in a P2P location-based service system, users may request hotel in-

formation within a geographical area around a conference site, which can be modeled as a range query over the corresponding location. Range query processing may also be applied over single values as an advanced functionality of existing systems. For instance, in a music file sharing system, song files are identified through title and release year (e.g., {"Pink Floyd", "1982"}); a range query "search all Pink Floyd songs between year 1980 and 1990" will enhance the system with range-aware functionality.

Existing works on P2P range query processing typically rely on structured overlay network protocols (e.g., BATON [11] and P-Tree [5]). While these protocols guarantee range query processing to complete within a bounded number of routing hops, they are not widely adopted in practical systems due to the following drawbacks: (1) the construction and maintenance, particularly during network churn, of structured overlay networks are non-trivial; and (2) these approaches often disregard the potential of data caching during query processing, which may affect query execution performance significantly.

In contrast, unstructured P2P overlay network architecture is widely adopted by practical systems, where constrained flooding mechanisms are usually employed for keyword-based searching (e.g., in Gnutella and Bubblestorm [19]). Since data are cached at peers after retrieval, subsequent queries for the same data can be answered by multiple caches, facilitating the search process. In the case of range query processing, the queries can also be shipped within the network through flooding. Data retrieved after the query execution are easily cached at query issuers, which produce *distributed range caches* that can potentially be used during subsequent range query processing.

Since unstructured P2P overlay networks are built without the knowledge of data placement, average communication cost (e.g., the number of messages or latency) per query is inversely proportional to the number of data copies, or replicas, in the network that satisfy the query [4]. Those data results that are not well-replicated may incur a higher cost to retrieve, delaying the progress of the range query processing.

It is often the case that range query results include

¹<http://www.gnutella.com>

²<http://www.bittorrent.com/>

³For conciseness, we use *range query* and *range selection query* interchangeably in the remainder of the paper.

⁴<http://www.joost.com>

data items that are not well replicated. For example, in a location-based hotel reservation system in P2P networks, “popular” hotels, such as those close to a conference site, are usually queried first by conference participants. When these hotels are booked in full, users tend to relax the range constraints (e.g., with respect to geographical proximity) to explore other hotel information, which may not be well replicated in the network yet. For simplicity, we refer to these data items that are not well replicated as *poorly replicated* data items. The approach proposed in this work will predict and prefetch those poorly-replicated data items that may potentially be requested in subsequent range queries and then facilitate the caching of these data to improve overall query execution performance.

Existing approaches are not effective in facilitating the caching of the poorly-replicated data items that may potentially be queried in the future. In unstructured P2P overlay networks, *uniform replication scheme* (deployed in KaZaa) caches each data item at a fixed number of peers so that it preventively excludes poorly-replicated data items from the system. However, this scheme is oblivious to the knowledge of query distribution such that the fixed number of peers that manage the data items covered by “popular” queries may be overloaded. Moreover, the replication scheme requires strong altruistic cooperation from peers in that peers may cache data items regardless of whether they are issuing queries, which may not be feasible in uncooperative P2P environments. In contrast, with the *proportional replication scheme* that is employed in Gnutella⁵, each peer only caches results after query execution such that caching process is triggered by query processing and only query issuers rather than arbitrary peers cache the results. Similarly, *square-root replication scheme* [4] makes the number of cached data items proportional to the square root of the number of the corresponding queries, improving average query processing performance with respect to constrained cache sizes. In comparison to the uniform replication scheme, the last two schemes achieve better load-balancing by considering query distributions and do not rely on altruism from peers. However, they disregard the caching of poorly-replicated data items, which may bottleneck the performance of the range query since these data are part of the range query result.

The key to the problem is to recognize those poorly-replicated data items that may potentially be requested by range queries in the future, and to facilitate the caching of these data items. With respect to unstructured P2P overlay networks, the following design principles are crucial: (1) the approach needs to be purely decentralized; (2) the approach should be efficient, without incurring significant communication or computational overhead; and (3) it is desirable that the approach can be deployed with existing routing protocols and replication schemes that have been devel-

oped for unstructured P2P networks.

Using the intuition that well-replicated data are “popular” data cached in the P2P system, we propose an efficient, decentralized *popularity-aware prefetch-based* approach to facilitate the caching of poorly-replicated data items. Prefetching in the context of other application domains (e.g., compiler technology) has been well-studied. Our approach, in the context of prefetching range data in P2P systems, is novel in that it is data-popularity aware: (1) peers independently collect global information about the relationship between poorly-replicated data items and “popular”, well-replicated, ones involved in previously executed queries, enabling an adaptive approach for range query processing (see Section 3 for details); (2) since popular data are easily obtained from peers through flooding or random walk mechanism, simply prefetching poorly-replicated data items can be more cost-effective with respect to bandwidth consumption; and (3) sufficient query issuers will exist over “popular” data items, providing opportunities to piggyback “correlated” poorly-replicated data items onto popular data and thereby facilitating the prefetching process. Briefly, the contributions of this work include the following.

- This is the first work that addresses distributed range caching problems in unstructured P2P overlay networks. In particular, purely decentralized mechanisms are developed to locate poorly-replicated, “correlated”⁶, data items that are potentially queried in the future, and an adaptive prefetch-based approach is proposed to improve performance of the range query processing that involves those poorly-replicated data items.
- The effectiveness of the approach is demonstrated theoretically by proving that, under a specific query distribution model with an increasing query load, the approach can improve overall performance of the queries that retrieve poorly-replicated data items by at least a factor of $\mathcal{O}(\ln m)$, where m is the number of queries, even when network churn and cache expiration exist.
- Through extensive simulations, the effectiveness of the approach is demonstrated under various query load settings.

The organization of the remainder of the paper is as follows. In Section 2 we analyze performance of cache-based range query processing in unstructured P2P overlay networks. Section 3 covers the prefetch-based approach that results in substantial communication cost savings for range query processing. Performance evaluation results are presented in Section 4. We review the related work in Section 5 and conclude this paper in Section 6.

2 Cache-based Range Query Processing

Range queries typically involve the range constraints that are defined over numeric-valued data [14]. The processing

⁵<http://www.gnutella.com>

⁶We define “correlation” in Section 3.1.

of a specific range query completes when all distinct result data items (either from original data sources or from caches) that satisfy the range constraints have been retrieved. Range queries can be issued by any peer in the system, and be shipped to other peers through flooding, gossip, or random walk routing mechanisms. Since flooding is no more effective than random walk with respect to routing cost [4], while gossip mechanism is based on random walk in that queries are shipped to randomly chosen peers in the network, we focus on the random walk mechanism in the remainder of this paper without loss of generality.

Based on the random walk mechanism, both the number of messages and the latency to retrieve a specific data item are inversely proportional to the number of the data item replicas in the network. Suppose that the cost function (denoted by $Cost$) of range query processing is defined over $Q \times R$, where Q denotes the set of range queries and R is the real value domain of query processing cost⁷. Given a query $q \in Q$, $Cost(q) \propto \frac{1}{r_q}$ (i.e., $Cost(q)$ is proportional to $\frac{1}{r_q}$), where r_q is the lower-bound of the number of data item replicas that can satisfy q . When there exist multiple distinct data items $\{s_1, s_2, \dots\}$ in the network that satisfy range query q , $r_q = \min(|s_1|, |s_2|, \dots, |s_i|, \dots)$, where we denote by $|s_i|$ the corresponding numbers of data item s_i replicas.

We are especially concerned about the range queries that include poorly-replicated data items because the overall range query performance is affected by them. We refer to the period that specific data items have not been sufficiently replicated as the *cold period*. Suppose that the results of a range query q include data item s , which initially has a single replica (i.e., the original data item itself) in the network; when there are m subsequent q queries issued, the overall query execution cost during s ' cold period is computed as below, based on the well-established proportional and square-root replication schemes respectively. For simplicity, we suppose that each execution of the query is initiated by a distinct peer.

- *Proportional replication scheme* Each query execution is expected to increase the number of replicas by one, such that the overall query processing cost with respect to m q queries, denoted by $Cost(q)$, is derived as below. Because sequence (3) does not converge, we present an approximate result with respect to a considerably large m .

$$Cost(q) = \sum_{i=1}^m \left(\frac{N}{i}\right) \quad (1)$$

$$= N + \frac{N}{2} + \frac{N}{3} + \dots + \frac{N}{m} \quad (2)$$

$$\approx N \times \ln m \quad (3)$$

⁷In this work, we focus on the query shipping cost to locate query results, ignoring local processing cost.

- *Square-root replication scheme* Although the square-root replication scheme performs better than the proportional replication scheme with respect to range query processing performance under constrained overall cache sizes [4], the range query processing involving poorly-replicated data items may incur higher communication cost. The following derivation presents the overall query processing cost, where sequence (5) increases monotonically and never converges. For simplicity, we assume that the number of replicas exactly equals the square root of the corresponding number of queries⁸, which does not affect the validity of the obtained result.

$$Cost(q) = \sum_{i=1}^m \left(\frac{N}{\sqrt{i}}\right) \quad (4)$$

$$= N + \frac{N}{\sqrt{2}} + \frac{N}{\sqrt{3}} + \dots + \frac{N}{\sqrt{m}} \quad (5)$$

The above analysis shows that range query processing performance is affected when queries retrieve poorly-replicated data items during their cold period. In the next section, we propose the prefetch-based approach, where poorly-replicated data items are prefetched by query issuers that request the well-replicated data items ‘‘correlated’’ to the poorly-replicated ones. This potentially decreases the retrieval cost of the poorly-replicated data items within cold periods and improves the range query execution performance.

3 Prefetch-based Caching

In this section, we first define data *correlation*, which materializes the locality concept that is essential to prefetch-based mechanisms [18]. Then we detail the design of the prefetch-based approach.

3.1 Data Correlation

To quantify the correlation between poorly-replicated data items and well-replicated ones, we introduce a distance function D . For Euclidean range space, when data items represent point data (e.g., the longitude and latitude information of locations), D can simply be the one-dimensional Euclidean distance function between the point values⁹. Instead, when data items correspond to range segments (e.g., the range span of longitude and latitude information of a region around a specific location), D may be defined over the Euclidean distance between the centroid points (e.g., median points in one-dimensional space) of the corresponding range segments. Other applications may employ their customized distance functions, which does not affect the applicability of the prefetch-based approach.

Based on the distance function D , data *correlation* is defined as follows: *two data items s and s' are correlated*

⁸The actual number of replicas equals the square root of the corresponding query load size multiplied by a constant factor [4].

⁹This does not conflict with the focus on range query processing since range queries may include multiple point values.

if $D(s, s') \leq \tau$. The correlation threshold τ can be pre-defined and configured by peers when they join the overlay network, which may not be sufficiently flexible since the threshold may be over or under-valued. For example, with respect to a specific range query load, when τ is set too low, data items that are covered by the same range queries may potentially be regarded uncorrelated; in contrast, when τ is set too high, more irrelevant data items may be regarded correlated even if they are never queried together. Since we are especially interested in the piggybacking of poorly-replicated data items with the well-replicated query results, τ is measured based on the distances between poorly-replicated data items and well-replicated ones from the history of executed range queries.

On one hand, this approach takes the information of range queries into account such that: (1) the inherent correlation of data items within the same range query is captured; and (2) since intuitively the well-replicated data items are “popular” among peers, their correlated (poorly-replicated) data items may also become “popular” with a higher probability, which is recognized by our approach. On the other hand, the approach is popularity-aware, enabling *adaptive* data prefetching: the greater the portion of range query workload that retrieves poorly-replicated data items, the more precisely τ captures the expected distances between poorly-replicated data items and well-replicated ones. Conversely, when range queries seldom retrieve poorly-replicated data items, τ tends to be close to zero such that prefetching may not even be triggered.

While there do exist correlations between poorly-replicated data items, or between well-replicated data items, these are not considered in this work because: (1) the proposed approach relies on query issuers requesting well-replicated data items to prefetch poorly-replicated ones; the correlation between poorly-replicated data items is less important because it does not indicate the involved data items will be requested in subsequent queries; and (2) the retrieval of well-replicated data items incurs low processing cost, such that the prefetching of these data items may not be cost-effective with respect to bandwidth consumption. Performance evaluation (Section 4) supports the belief that simply prefetching all correlated data items regardless of data popularity may not be efficient.

Specifically, each peer records the history of the range queries that it issued previously. Each history record contains the maximum distance between any poorly-replicated data items and well-replicated ones that are involved in the range query at query execution time. Each peer then randomly samples a configurable number of other peers in the overlay network and obtains their query processing history records. This captures approximate global information about how poorly-replicated data items affect range query execution. The random sampling in unstructured

P2P overlay networks can be realized through existing techniques [7, 12]. Query processing history records are collected by peers periodically such that they learn up-to-date knowledge on executed range queries.

Once a set of distance values from the history records, denoted by $\{d_i\}$, are obtained, each peer p estimates the distribution of the distance values through kernel estimation technique [17], which is a non-parametric data distribution modeling scheme. Non-parametric modeling schemes are advantageous in P2P environments since no a-priori knowledge about data distribution is required. Then we set τ to be the expected distance value under the distribution model. In this work, the commonly-used Gaussian kernel function $K(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}x^2}$ is employed, and the estimated probability density function (PDF) is $PDF_k(x) = \frac{1}{Sh} \sum_{i=1}^S K(\frac{x-x_i}{h})$, where S denotes the number of history records, x_i denotes the distance value corresponding to the i_{th} sample, and h is a smoothing factor that is configurable with respect to specific applications. The derivation of τ is shown below, where $y_i = \frac{x-x_i}{h}$ for each specific i .

$$\begin{aligned} \tau &= \int_0^\infty x \times PDF_k(x) dx = \int_0^\infty x \times \frac{1}{Sh} \sum_{i=1}^S K(\frac{x-x_i}{h}) dx \\ &= \frac{1}{Sh\sqrt{2\pi}} \times \sum_{i=1}^S (\int_0^\infty x \times e^{-\frac{1}{2}(\frac{x-x_i}{h})^2} dx) \end{aligned}$$

Define $y_i = x - x_i$,

$$\begin{aligned} \tau &= \frac{1}{Sh\sqrt{2\pi}} \times \sum_{i=1}^S (\int_0^\infty x \times e^{-\frac{1}{2}y_i^2} \times h dy_i) \\ &= \frac{1}{S\sqrt{2\pi}} \times \sum_{i=1}^S (\int_0^\infty (hy_i + x_i) \times e^{-\frac{1}{2}y_i^2} dy_i) \\ &= \frac{1}{S\sqrt{2\pi}} \times \sum_{i=1}^S (h \int_0^\infty hy_i \times e^{-\frac{1}{2}y_i^2} dy_i + \int_0^\infty x_i \times e^{-\frac{1}{2}y_i^2} dy_i) \\ &= \frac{h}{\sqrt{2\pi}} + \frac{1}{2S} \times \sum_{i=1}^S x_i \end{aligned}$$

With respect to range query processing, since query issuers and peers holding query results may obtain different values of τ , we always choose the query issuer’s τ as the correlation threshold, because they can flexibly adjust the value of τ to include other constraints (*e.g.*, local storage constraint).

3.2 Data Popularity

The computation of τ requires the knowledge of data *popularity*, which is used to distinguish poorly-replicated data items from well-replicated ones. Although the obtained query history records can be employed to estimate data popularity [23], they may not be sufficiently accurate to reflect the overall data distribution in the overlay network. Thus, we consider a purely decentralized approach to estimate data popularity directly. Each peer evaluates the popularity of the local data through an exploration process. Consider a peer p with a set of data items, denoted by $\mathcal{S} = \{s_1, s_2, \dots\}$. p issues an “exploration query” over each $s_i \in \mathcal{S}$, where a configurable Time-to-Live (TTL) counter is attached to each exploration query. During each hop, the

TTL counter decreases by one and all appearances of s_i are recorded by p . The shipping of an exploration query terminates when TTL equals zero. Once the exploration process completes for all data items, p figures out the number of cached replicas for each s_i during the exploration process. Those s_i with more replicas (*i.e.*, above threshold T) are regarded as “well-replicated” and the others are considered “poorly-replicated”. A similar approach has been employed in PIER to decide data popularity [9]. Both the TTL and T are configurable, where the TTL value is usually set to be small to avoid large explorations.

3.3 Prefetch-based Approach

Suppose that peer p receives a range query issued by p' that covers a data item s ; p will reply to p' with all its cached data items $\{s_1, s_2, \dots, s_i, \dots\}$ that satisfy the following conditions: (1) each s_i is correlated to s based on the distance function D and the correlation threshold τ that is attached to the query (*i.e.*, $D(s, s_i) \leq \tau$); and (2) based on the data popularity measurement of p' (*i.e.*, the query issuer), s is well-replicated and s_i is poorly-replicated.

For instance, consider multiple query issuers requesting data items s and s' through queries q_1 and q_2 respectively, as illustrated by white and grey nodes in Figure 1(a). For simplicity, suppose that s and s' are cached at peer p that receives the queries. Then Figure 1(b) demonstrates the proportional replication scheme and Figure 1(c) shows the proportional replication scheme enhanced with prefetching, where peers denoted by dark nodes eventually cache the poorly-replicated data item s' . Due to the prefetching of s' by query issuers over q_1 , s' is cached more quickly in Figure 1(c).

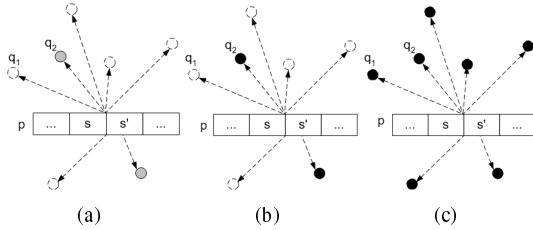


Figure 1. Caching approaches

We describe the prefetch-based approach in Algorithm 1 regarding peer p when it receives a query q covering data item s that is issued at p' . The approach has several advantages: (1) it makes slight changes to existing replication schemes without adjusting the overlay network architecture or routing algorithms; all decisions, including the computation of the data correlation threshold and data popularity, are made independently without costly coordination among peers; (2) poorly-replicated data items are piggybacked during the query processing over well-replicated data items, saving communication cost; moreover, peers are not required to be altruistic since only query issuers who utilize the query processing services take on the prefetching over-

head cost; and (3) the approach is purely decentralized without relying on any centralized mechanisms, thus providing scalability.

Algorithm 1 *prefetch_based_caching*(q, τ)

```

1: if  $p$  learns that  $s$  is well cached then
2:   for all data items  $s'$  that are poorly-replicated and  $D(s, s') \leq \tau$  do
3:     if  $D(s, s') \leq \tau$  then
4:        $p$  replies to  $p'$  the data item  $s$ ;
5:     end if
6:   end for
7: end if

```

3.4 Guarantees on Performance Improvement

With the prefetch-based approach, the population of poorly-replicated data items is affected by the query volume over the correlated well-replicated data items. We now show that, under the following query distribution, our approach guarantees that the overall range query cost over poorly-replicated data items is $N \times \mathcal{O}(1)$, which is at least $\mathcal{O}(\ln m)$ factor less than the counterpart when no prefetching is enforced.

Consider a range query q_1 involving data item s and query q_2 involving data item s' , where s is well-replicated while s' is poorly-replicated. Suppose that, initially only one replica of s' exists in the network, while there exist n replicas of s including original and cached ones. Consider a query load consisting of m_1 q_1 queries and m_2 q_2 queries. For simplicity, all q_1 (and q_2) queries are uniformly distributed across a certain period of time; thus it is expected that $\frac{m_1}{m_2}$ q_1 queries are processed when each q_2 is processed. Note that this assumption is not necessary for the following analysis; it will be clear shortly that only if the number of q_1 queries is sufficiently large compared with that of q_2 queries, the analysis will hold. We also assume that the number (n) of s replicas is stable, and will discuss the case when n may change. We now prove Theorem 1, where $\delta > 1$ is a system parameter that affects the constant factor of the average query operation cost (*i.e.*, $\mathcal{O}(1)$), and N denotes the network size. A proof is presented subsequently.

Theorem 1. When $\frac{m_1}{m_2 \times n} \geq 2^\delta$, the overall cost for processing m_2 q_2 queries equals $\mathcal{O}(1) \times N$.

Proof. Suppose that initially, there is one replica of s' and it is cached together with s on a peer. This does not affect the generality of the analysis since subsequently prefetched s' will be cached together with s . During each period, one q_2 is processed and $\frac{m_1}{m_2}$ q_1 queries are issued such that the expected number of q_1 queries that visit a peer holding s' equals $\frac{m_1}{m_2 \times n}$. When $\frac{m_1}{m_2 \times n} \geq 2^\delta$, the number of cached s' data item copies increases by 2^δ fold. Recursively, suppose $i > 1$, when there are $(i-1)^\delta$ s' replicas in the network after the execution of the $(i-1)_{st}$ q_2 query, the expected number of peers that cache s' through the execution of query q_1 equals $\frac{m_1}{m_2 \times n} \times (i-1)^\delta$ and is no less than i^δ , as shown in the following derivation.

$$\frac{m_1}{m_2 \times n} \times (i-1)^\delta > (2 \times (i-1))^\delta \quad (6)$$

$$\geq \left(1 + \frac{1}{i-1}\right) \times (i-1)^\delta = i^\delta \quad (7)$$

Consequently, when the execution of the i_{th} q_2 query completes, no less than i^δ peers will cache s' . The overall shipping cost to resolve all q_2 queries is then computed as follows. While q_2 queries are executed sequentially, the above analysis is easily extended to handle concurrent query execution.

$$Cost(q_2) \leq \sum_{i=1}^{m_2} \frac{N}{i^\delta} \quad (8)$$

$$= N + \frac{N}{2^\delta} + \frac{N}{3^\delta} + \dots + \frac{N}{m_2^\delta} \quad (9)$$

Since sequence (9) is a Riemann-Zeta sequence [6], it converges for all real values $\delta > 1$. Thus the overall query operation performance with respect to the query load of q_2 is $Cost(q_2) = \mathcal{O}(1) \times N$ when m_2 is considerably large. \square

For example, when $\delta = 1.5$, $Cost(q_2) \approx 2.612 \times N$. This overall cost is at least $\mathcal{O}(\ln m)$ times less than those achieved by existing replication schemes (addressed in Section 2) that do not employ prefetching.

A sufficient condition of Theorem 1 is that the ratio of the number of q_1 queries over the number of s caches (denoted by n) is no less than 2^δ . However, multiple factors may affect this number in practice: (1) after the execution of q_1 queries, peers are capable of caching s , increasing n ; (2) under network churn, peers may fail (or leave) suddenly, decreasing n ; and (3) when cache expiry schemes are deployed in P2P systems for data freshness [1, 10], n may also change.

Cache replacement may affect the value of n as well. Specifically, if the cache size is sufficiently large, peers can hold all recently cached data items and cache replacement would happen infrequently. If the cache size is limited, cache replacement behavior can be easily integrated with cache expiry. For example, a widely employed recency-based cache replacement strategy such as Least-Recently-Used can be enforced through the cache expiry process by choosing an appropriate expiry period. This would evict least-recently-used data items from caches after each expiry period.

Next, given a specific query load of q_1 , Theorem 1 can be satisfied even when network churn and/or cache expiry occur. We prove this in Theorem 2. Recall that, during each period, one q_2 query is issued. Any peer (and cache) fails (or leaves) with a probability of $0 \leq r \leq 1$ per period, and all cached data items expire after $k > 0$ periods of time. $l(i)$ represents the expected number of q_1 queries during the i_{th} period, $n(i)$ denotes the number of s data item replicas after the i_{th} period, and $n(0) = n$ denotes the initial number of s replicas.

Theorem 2. *When the query load of q_1 satisfies the following distribution,*

$$l(i) = \begin{cases} 2^{\delta+1}(1+2^{\delta+1})^{i-1}(1-r)^{i-1}n, & \text{when } i < k \\ 2^{\delta+1}(1+2^{\delta+1})^{i-1}(1-r)^{i-1}n - \mathcal{O}((1+2^{\delta+1})^{i-k}(1-r)^i), & \text{when } i \geq k \end{cases}$$

the number of s' replicas will be no less than i^δ after the i_{th} period, where $\delta > 1$; consequently, the overall cost of processing m_2 q_2 queries is bounded by $\mathcal{O}(1) \times N$, even under network churn and cache expiry.

Proof. For brevity, we only prove the case when $i \geq k$. It is easy to apply the same derivation for the case when $i < k$, where no cached data expires.

When $i \geq k$, the number of s replicas (denoted by $n(i)$) after the i_{th} period is computed below, where the first component consists of the accumulated number of s until the $(i-1)_{st}$ period, added by the increase of the s (denoted by $l(i)$) during this period; the second component covers the loss of s replicas due to cache expiry. Both components are multiplied by corresponding damping factors to reflect network churn.

$$\begin{cases} n(i) = (n(i-1) + l(i)) \times (1-r), & \text{when } i < k \\ n(i) = (n(i-1) + l(i)) \times (1-r) - l(i-k) \times (1-r)^k, & \text{when } i \geq k \end{cases}$$

Consider a more relaxed function $n'(i) = (n'(i-1) + l(i)) \times (1-r) - (1-r)^k$, where $n'(0) = n(0) = n$. Since $(1-r)^k \leq l(i-k) \times (1-r)^k$, it is obvious that $n'(i) \geq n(i)$. Then similarly consider,

$$\begin{cases} n'(i) = (n'(i-1) + l(i)) \times (1-r), & \text{when } i < k \\ n'(i) = (n'(i-1) + l(i)) \times (1-r) - (1-r)^k, & \text{when } i \geq k \end{cases}$$

Suppose $A = 2^{\delta+1}$ and make $l(i) = A \times n'(i-1)$. The following derivation computes $n'(i-1)$.

$$\begin{aligned} n'(i-1) &= (n'(i-2) + l(i-1)) \times (1-r) - (1-r)^k \\ &= (1+A)^{i-1} \times n'(0) \times (1-r)^{i-1} - \sum_{x=1}^{i-k-1} ((1+A)^x (1-r)^{x+k}) \\ &= (1+A)^{i-1} \times n \times (1-r)^{i-1} - \sum_{x=1}^{i-k-1} ((1+A)^x (1-r)^{x+k}) \\ &= [(1+A)(1-r)]^{i-1} \times n - \frac{(1+A)^{i-k} (1-r)^i}{(1+A)(1-r)-1} + \frac{(1-r)^k}{(1+A)(1-r)-1} \end{aligned}$$

With reasonable network churn rate, $(1+A)(1-r)$ will be larger than 1. By ignoring the last component $\frac{(1-r)^k}{(1+A)(1-r)-1}$, which is a small constant when k is fixed, $n'(i-1)$ is bounded by $[(1+A)(1-r)]^{i-1} \times n - \mathcal{O}((1+A)^{i-k}(1-r)^i)$. Consequently, $l(i) = A \times n'(i-1) = 2^{\delta+1} \times [(1+A)(1-r)]^{i-1} \times n - \mathcal{O}((1+A)^{i-k}(1-r)^i)$ is a sufficient condition for $l(i) \geq A \times n(i-1)$ because $n'(i-1) > n(i-1)$. Then the following analysis about the query processing cost holds based on mathematic recursion, where $i > 1$.

$$l(i) \times (i-1)^\delta \geq 2^{\delta+1} \times (i-1)^\delta \times n(i-1) \quad (10)$$

$$\geq 2 \times \left(\frac{i}{i-1}\right)^\delta \times (i-1)^\delta \times n(i-1) \quad (11)$$

$$\geq (i^\delta + (i-k)^\delta) \times n(i-1) \quad (12)$$

$$\rightarrow \frac{l(i) \times (i-1)^\delta \times (1-r)}{n(i-1) \times (1-r)} - (i-k)^\delta \geq i^\delta \quad (13)$$

Equation 13 shows that under the setting that the network churn rate equals r and the cached data items are evicted after k periods, the number of s' replicas (including the original and cached ones) is no less than i^δ after the i_{th} period. It is direct that, the overall cost for processing m_2 q_2 queries is no more than $\sum_{i=1}^{m_2} \frac{N}{i^\delta} = N + \frac{N}{2^\delta} + \frac{N}{3^\delta} + \dots + \frac{N}{m_2^\delta} = \mathcal{O}(1) \times N$, even under network churn and cache expiry. \square

4 Performance Evaluation

4.1 Experimental Methodology

To study range query processing performance under the prefetch-based approach, we use the discrete event simulator p2psim¹⁰. We generate a network topology by randomly mapping peers to coordinates in a square area using a uniform distribution (which also generates local skew effects in the mapping). We consider a network of up to $N = 3000$ peers. Both a static network setting and a dynamic setting with network churn are simulated.

We consider one-dimensional range queries within the domain of $R = [0, 10000]$. Initially each peer holds $i = 10$ items that are randomly chosen within the range. The range query load is synthesized as follows. The number of range queries within each *epoch*¹¹ follows Poisson distribution with mean $m = 100$; the query execution process runs for 20 epochs. To simulate the Power-law characteristics of query load that is common in real world [15], the range domain is evenly partitioned into a configurable number (*i.e.*, 10) of range segments; the start point (*i.e.*, lower bound) of each range query is generated over a randomly chosen segment based on the principles of *growth* and *preferential attachment*, which produce Power-law query distribution [2]. When a query is issued more than a number of times and becomes “popular”, peers issuing it start to migrate the query.

To cover various peer behaviors, three migration patterns are considered (in Figure 2), producing *transitional queries*, *relaxed queries* and *consecutive queries*. For each original query Q_p , a configurable number of migrated queries Q_r are generated and executed subsequently. The number (m) of the original queries and the number (m') of migrated queries per original query will be configured shortly.

- *Transitional queries*. All points in the data domain are potentially chosen as the start point of Q_r with a probability proportional to the Euclidean distance away from the start point of Q_p . The range span of Q_r follows Poisson distribution with mean $span = 100$. This migration pattern may simulate the scenario that, after finding that the hotel rooms around a specific city are all booked, peers may turn to check other proximate cities.
- *Relaxed queries*. The start point of each Q_r equals to that of Q_p , denoted by *start*. Then the range span

of Q_r is relaxed to follow Power-law distribution between $(span, |D| - start)$. This migration pattern covers those relaxed queries over different radius of the region around a specific location (*e.g.*, a conference site).

- *Consecutive queries*. The start point of each Q_r adopts the end point of Q_p . The range span of Q_r follows Poisson distribution with mean $span = 100$. This migration pattern captures the setting when peers adjust the search scope by retrieving the data that are disjoint but contiguous to initial geographical range constraints.

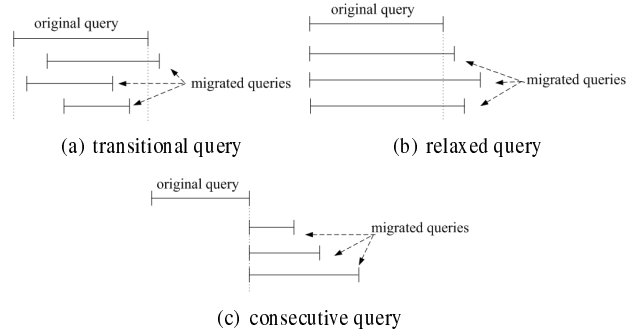


Figure 2. Synthesized Migrated Queries

Since the random walk mechanism requires a randomization mechanism, all experimental results are averaged over three runs, each with a different random seed. Recall that both the number of messages and the query routing latency are inversely proportional to the number of query result replicas in either flooding or random walks; we only report the number of messages consumed. Since only the reply message sent back to query issuers carry query results and prefetched data while all other messages only carry the query itself (*e.g.*, range constraints), the amortized message size is small. Thus, message size is not studied in the performance evaluation. In this simulation, we focus on the evaluation of the proportional replication scheme, while leaving that of the square-root replication scheme to future work.

4.2 Evaluation of Query Shipping

For simplicity, this experiment assumes a static environment, where no network churn or data insertion occur. We will consider network churn shortly in Section 4.4. Initially, $m = 100$ queries are generated during each epoch by following Power-law distribution. When a query becomes “popular”, $m' = 5$ migrated queries (respectively transitional, relaxed and consecutive queries) are generated based on each original query.

The average number of messages per query is shown in Figure 3. In addition to the *nonprefetch* option that acts the same as the proportional replication scheme and the *prefetch* option that corresponds to our prefetch-based approach, we also consider *allprefetch* approach, which prefetches not only correlated poorly-replicated data items but well-replicated ones during the execution of range

¹⁰p2psim: <http://pdos.csail.mit.edu/p2psim/>

¹¹In this experiment, each epoch lasts 5×10^6 milliseconds.

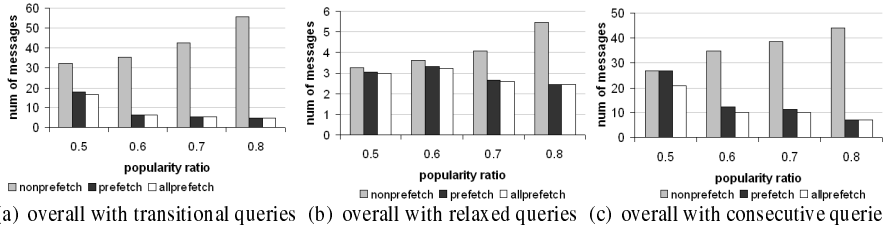


Figure 3. Average shipping cost per query

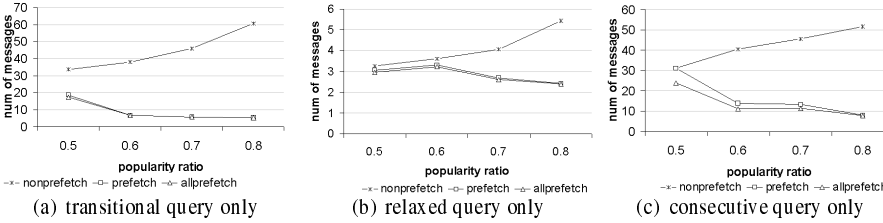


Figure 4. Average shipping cost per query

queries. For presentation, the results are illustrated with a logarithmic scale with base of 2.

The experimental results show that the prefetch-based approach effectively decreases the overall number of messages per query with respect to all migration patterns. In this experiment, the ratio (denoted by pr) corresponding to the x -axis of Figure 3 may trigger the query migration: when an original query is issued by peers for a sufficiently large number of times (*i.e.*, $c \times pr \times N$), migrated queries are generated for this query, where N denotes the network size and $c = 0.6$ for this simulation. When pr is higher, the number of original queries is expected to increase, such that more migrated queries are generated, potentially including more poorly-replicated data items; consequently the average query execution performance of the *nonprefetch* option goes up. In contrast, the average query execution cost under the *prefetch* option is lower under all settings, showing the effectiveness of our approach. In Figure 4(a)-(c) we also present the shipping cost consumed by only migrated queries, indicating that the performance improvement of our approach is primarily due to the cost savings over migrated queries. These figures also illustrate that the benefit of simply prefetching all correlated data items including well-replicated ones (under the *allprefetch* option) is not significant. We measure the bandwidth consumption of the *prefetch* option versus that of the *allprefetch* option, which can be 60% more. This confirms our belief that data popularity should be considered in prefetching.

Since “relaxed queries” are generated by following Power-law distribution and a small set of migrated queries are issued with an exponentially higher probability, the (initially) poorly-replicated data items covered by correspond-

m'	1	5	10
<i>transitional</i>	10.47	5.39	4.61
<i>relaxed</i>	2.66	2.31	2.30
<i>consecutive</i>	9.93	7.59	7.59

Table 1. Query execution cost with different query load

cache size	500	1000	1500	2000
<i>transitional</i>	357	321	255	241
<i>relaxed</i>	2.73	2.11	2.07	2.07
<i>consecutive</i>	72	24	16.8	12.19

Table 2. Query execution cost under various cache sizes

ing query results become well replicated very quickly. Thus the number of messages per query is significantly lower than that of the queries generated under the other two migrated patterns. This holds for other experimental results over “relaxed queries” in the remainder of the paper.

With respect to the load of migrated queries, we evaluate different numbers of migrated queries (*i.e.*, $m' = 1, 5$ and 10), with $pr = 0.7$. As shown in Table 1, the average number of messages decreases because a relatively larger number of queries (over poorly-replicated data items) may benefit from the prefetching of poorly-replicated data.

4.3 Evaluation with Cache Constraints

In practical systems, peers usually have caches with limited storage sizes. Moreover, cached data items may expire after a period of time, as discussed in Section 3. In this simulation, we measure the query shipping cost with respect to these cache constraints.

We consider 500, 1000, 1500 and 2000 cache entries (*i.e.*, the cached data items per peer). The experimental results are shown in Table 2, which indicate that when the cache size is larger, the average query shipping cost tends to be lower. This is primarily due to the fact that larger caches can hold more prefetched data, facilitating query execution. We also evaluate the query processing performance under various cache expiry periods (*i.e.*, $5 \times 10^6, 1 \times 10^7, 5 \times 10^7$ and 1×10^8 simulation milliseconds). The results (shown in Figure 5) demonstrate that when the lifespan of cached data items increases, the average number of messages per query decreases. This is because longer cache life means that data expire less frequently, leading to better use of cached data items.

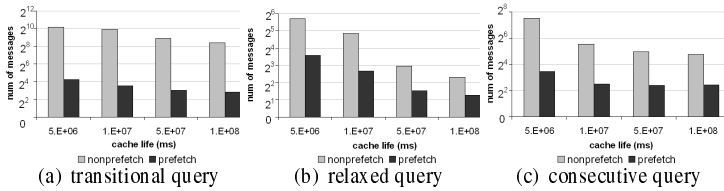


Figure 5. Query execution cost per query under cache expiry

	m'	
	1	10
<i>transitional (byte)</i>	87500	72981
<i>relaxed (byte)</i>	18	2
<i>consecutive (byte)</i>	40652	38274

Table 3. Bandwidth consumption per query with different migration queries

4.4 Evaluation under Dynamic Settings

In P2P networks, peers may leave and fail arbitrarily, affecting the number of cached data items in the network. Moreover, peers may join the network with new data items and data insertion manipulation may also be issued by existing peers.

We simulate the network churn by allowing each peer to leave with a probability of r during each epoch, where r varies from 0.1 to 0.4. To make the network stable, peers join the network at the same rate. The number of messages per query is shown in Figure 6. The experimental results indicate that when failure rate increases, the number of messages per query goes up, which is due to more random walks being required for completing the query processing during network churn.

To measure range query processing performance when data insertion occurs, peers generate a configurable number (*i.e.*, 1, 10, 100 and 1000 per epoch) of new data items that are randomly chosen within the data domain. In this experiment, no migrated queries are generated. As shown in Figure 7, the average shipping cost per query increases steadily during the insertion of new data items. This is not surprising, because when peers keep issuing the same range queries, they still need to locate the new data items that initially are poorly-replicated.

4.5 Evaluation of Bandwidth and Adaptivity

The correlation threshold τ is computed adaptively based on the knowledge of executed range queries. In our performance studies, we compare the average bandwidth costs of data prefetching per query with respect to a query load consisting of the same number of queries but with different migration query loads: either $m' = 1$ migrated query is generated for a “popular” original query or $m' = 10$ migrated queries are generated. The results shown in Table 3 indicate that the bandwidth consumed by the prefetch-based approach with a relatively larger number of migrated queries is higher. This is because the execution of more queries is affected by poorly-replicated data items and τ tends to be larger, increasing the volume of prefetched data. In contrast, when migrated queries are fewer, the value of τ is smaller. Thus the volume of prefetched data is lower, demonstrating the adaptivity of our approach.

5 Related work

Various architectures have been proposed to support large-scale P2P applications, including unstructured, super-

peer-based, and structured architectures [20]. In this work, we focus on unstructured P2P architecture, which imposes little constraints on the overlay network structure and has been widely employed in practical systems such as Gnutella, KaZaa, and BitTorrent.

Caching techniques have been employed in P2P systems [21], focusing on the routing efficiency rather than the range query processing that is addressed in this work. Distributed range caching mechanisms have also been developed for P2P networks [13, 16]. Kothari *et al* propose a distributed tree-based index to manage all range caches to facilitate the search of range data in P2P networks [13]. Sahin *et al* employ a hyper-rectangle-based overlay network to index multi-dimensional range data to support efficient range query processing [16]. These works use structured overlay networks in indexing range caches, which may potentially be affected by their non-trivial construction and maintenance cost. In contrast, our approach is focused on range caching in unstructured P2P overlay networks.

In unstructured P2P networks, effective search strategies are developed based on constrained flooding [22] and random walks [7]. A data replication scheme is employed in Bubblestorm [19], which does not take query distribution into consideration, potentially leading to load-balancing problems. In contrast, proportional and square-root replication schemes [4] are devised that consider query distribution. In this work, we consider a prefetch-based approach to facilitate the caching of poorly-replicated data items to improve range query processing performance, which can be directly deployed with these replication schemes. Moreover, our approach may prefetch different volumes of data items according to the knowledge of executed range queries, enhancing the adaptivity of the approach. An adaptive replication scheme has been studied in P2P networks [8], which however is focused on employing server load measurements to reduce replication cost.

Prefetching has been used in operating systems to facilitate instruction feed to CPU by fetching all possible instructions related to a branch conditional test in advance [18]. Our approach handles P2P range query processing and exploits data popularity to improve effectiveness of the prefetching. P2P media streaming systems also employ prefetching techniques to buffer upcoming stream data for smooth playback of the stream [3], which focus on adjust-

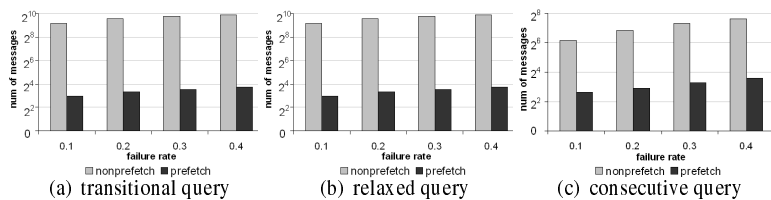


Figure 6. Query execution cost per query under network churn

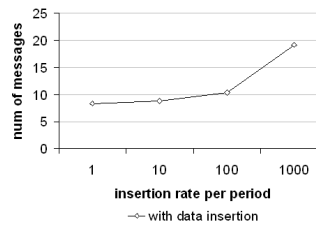


Figure 7. Query execution cost per query with data insertion

ment of the volume of prefetched data and are specific to streaming systems.

6 Conclusion

Under unstructured P2P overlay network architecture, query processing performance is usually decided by the number of query result replicas in the network. When range queries involve poorly-replicated data items, query execution performance degrades. In this work, we propose a popularity-aware prefetch-based caching approach that effectively facilitates the caching of poorly-replicated data items that are correlated with well-replicated ones, resulting in cost savings for future queries that access the poorly-replicated data. Our approach does not require strong altruistic cooperation from peers since only query issuers that use the query processing services incur prefetching overhead. Under various query load settings, we prove that the performance of range queries involving poorly-replicated data is guaranteed to improve. Experimentally, we also show that our proposed prefetch-based approach delivers substantial query processing cost savings.

References

- [1] W. Balke, W. Nejdl, W. Siberski, and U. Thaden. Progressive distributed top-k retrieval in peer-to-peer networks. In *Proc. Int. Conf. on Data Engineering*, pages 174–185, 2005.
- [2] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286:509–512, October 1999.
- [3] B. Cheng, X. Liu, Z. Zhang, and H. Jin. A measurement study of a peer-to-peer video-on-demand system. In *Peer-to-Peer Systems, First International Workshop*, 2007.
- [4] E. Cohen and S. Shenker. Replication strategies in unstructured peer-to-peer networks. In *Proc. ACM SIGCOMM*, pages 177–190, 2002.
- [5] A. Crainiceanu, P. Linga, J. Gehrke, and J. Shanmugasundaram. Querying peer-to-peer networks using P-Trees. In *Proc. 7th Int. Workshop on the World Wide Web and Databases (WebDB)*, pages 25–30, 2004.
- [6] H. M. Edwards. *Riemann’s Zeta Function*. Academic Press, 1974.
- [7] C. Gkantsidis, M. Mihail, and A. Saberi. Random walks in peer-to-peer networks. In *Proc. 23rd Annual Joint Conference of the IEEE Computer and Communications Societies*, 2004.
- [8] V. Gopalakrishnan, B. Silaghi, B. Bhattacharjee, and P. Keleher. Adaptive replication in peer-to-peer systems. In *Proc. 24th Int. Conf. on Distributed Computing Systems*, pages 360–369.
- [9] R. Huebsch, J. M. Hellerstein, N. Lanham, B. T. Loo, S. Shenker, and I. Stoica. Querying the internet with PIER. In *Proc. 29th Int. Conf. on Very Large Data Bases*, pages 321–332, 2003.
- [10] S. Iyer, A. I. T. Rowstron, and P. Druschel. Squirrel: a decentralized peer-to-peer web cache. In *Proc. ACM SIGACT-SIGOPS Symp. on Principles of Dist. Comp.*, pages 213–222, 2002.
- [11] H. V. Jagadish, B. C. Ooi, and Q. H. Vu. BATON: A balanced tree structure for peer-to-peer networks. In *Proc. 31th Int. Conf. on Very Large Data Bases*, 2005.
- [12] M. Jelasity, S. Voulgaris, R. Guerraoui, A.-M. Kermarrec, and M. van Steen. Gossip-based peer sampling. *ACM Trans. Comput. Syst.*, 25(3), 2007.
- [13] A. Kothari, D. Agrawal, A. Gupta, and S. Suri. Range addressable network: A P2P cache architecture for data ranges. In *Peer-to-Peer Computing*, pages 14–22, 2003.
- [14] R. Ramakrishnan and J. Gehrke. *Database Management Systems*. McGraw-Hill, 2002.
- [15] V. Ramasubramanian and E. G. Sirer. The design and implementation of a next generation name service for the internet. In *Proc. ACM SIGCOMM*, pages 331–342.
- [16] O. D. Sahin, A. Gupta, D. Agrawal, and A. E. Abbadi. A peer-to-peer framework for caching range queries. In *Proc. 20th Int. Conf. on Data Engineering*, pages 165–176, 2004.
- [17] D. Scott. *Multivariate Density Estimation: Theory, Practice and Visualization*. Wiley-Sons, 1992.
- [18] W. Stallings. *Operating Systems: Internals and Design Principles*. Prentice Hall, 2004.
- [19] W. W. Terpstra, J. Kangasharju, C. Leng, and A. P. Buchmann. Bubblestorm: resilient, probabilistic, and exhaustive peer-to-peer search. pages 49–60, 2007.
- [20] P. Valduriez and E. Pacitti. Data management in large-scale P2P systems. In *High Performance Computing for Computational Science - VECPAR 2004, 6th International Conference*, pages 104–118, 2004.
- [21] C. Wang, L. Xiao, Y. Liu, and P. Zheng. DiCAS: An efficient distributed caching mechanism for P2P systems. *IEEE Trans. Parallel Distrib. Syst.*, 17(10):1097–1109, 2006.
- [22] B. Yang and H. Garcia-Molina. Improving search in peer-to-peer networks. In *Proc. 22nd Int. Conf. on Distributed Computing Systems*, pages 5–12, 2002.
- [23] R. Zhang and Y. C. Hu. Assisted peer-to-peer search with partial indexing. In *The 24th Annual Joint Conference of the IEEE Computer and Communications Societies*, pages 1514–1525, 2005.