

2006

## Multi-Quality Data Replication in Multimedia Databases

Yi-Cheng Tu

Sunil Prabhakar

*Purdue University*, [sunil@cs.purdue.edu](mailto:sunil@cs.purdue.edu)

Jingfeng Yan

Gang Shen

Report Number:

06-009

---

Tu, Yi-Cheng; Prabhakar, Sunil; Yan, Jingfeng; and Shen, Gang, "Multi-Quality Data Replication in Multimedia Databases" (2006). *Department of Computer Science Technical Reports*. Paper 1652.  
<https://docs.lib.purdue.edu/cstech/1652>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.  
Please contact [epubs@purdue.edu](mailto:epubs@purdue.edu) for additional information.

**MULTI-QUALITY DATA REPLICATION  
IN MULTIMEDIA DATABASES**

**Yi-Cheng Tu  
Sunil Prabhakar  
Jingfeng Yan  
Gang Shen**

**CSD TR #06-009  
May 15, 2006**

# Multi-Quality Data Replication in Multimedia Databases

Yi-Cheng Tu, Sunil Prabhakar, Jingfeng Yan \* and Gang Shen †

## ABSTRACT

Quality is an essential property for multimedia databases. In contrast to other database applications, multimedia data can have a wide range of quality parameters such as spatial and temporal resolution, and compression format. Users can request data with a specific quality requirement due to the needs of their application, or the limitations of their resources. The database can support multiple qualities by converting data from the original (high) quality to another (lower) quality to support a user's query, or pre-compute and store multiple quality replicas of data items. On-the-fly conversion of multimedia data (such as video transcoding) is very CPU intensive and can limit the level of concurrent access supported by the database. Storing all possible replicas, on the other hand, requires unacceptable increases in storage requirements. Although replication has been well studied, to the best of our knowledge, the problem of multiple-quality replication has not been addressed. In this paper we address the problem of multiple-quality replica selection subject to an overall storage constraint.

We establish that the problem is NP-hard and provide heuristic solutions under two different system models: Hard-Quality, and Soft-Quality. Under the soft-quality model, users are willing to negotiate their quality needs, as opposed to the hard-quality system wherein users will only accept the exact quality requested. The hard-quality problem is reduced to a 0-1 Knapsack problem and we propose an efficient solution that minimizes the probability of request rejection due to unavailability of the requested quality replica. For the soft-quality system, an important optimization goal is to minimize utility loss. We propose a powerful greedy algorithm to solve this problem. Extensive simulations show that our algorithm performs significantly better than other heuristics. The algorithm is flexible in that it can be extended to deal with problems of distributed data replication

\*Y-C. Tu, S. Prabhakar, and J. Yan are with the Department of Computer Sciences, Purdue University, 250 N. University St. West Lafayette, IN 47907-2066, U.S.A. (emails: {tuyc, yanj, sunil}@cs.purdue.edu)

†G. Shen is with the Department of Statistics, Purdue University, 150 N. University St. West Lafayette, IN 47907-2067, U.S.A. (email: gshen@stat.purdue.edu)

and changes in query pattern.

## Keywords

Quality adaptation, integer programming, data replication, heuristic algorithm

## 1. INTRODUCTION

Quality is an essential property for multimedia databases. In contrast to other database applications, multimedia data can have a wide range of quality parameters such as spatial and temporal resolution, and compression format. Quality-aware multimedia systems [10, 23, 24, 26, 30] allow users to specify the *quality* of the media to be delivered based on their practical needs and resource availability on the client-side devices [23, 24]. The quality parameters of interest also differ by the type of media we deal with. For digital video, the quality parameters of interest include resolution, frame rate, color depth, signal-to-noise ratio (SNR), audio quality, compression format, and security level [30]. For example, a video editor may request a video at very high resolution when editing it on a high-powered desktop machine, but request the video at low resolution and frame rate when viewing it using a PDA. Different encoding formats may be desirable for different applications.

From the point of view of a video database, satisfying user quality specifications can be achieved using two complementary approaches<sup>1</sup>: i) store only the highest resolution copy, and convert it to the quality format requested by the user as needed at run-time; or ii) pre-compute each different quality that can be requested and store them on disk. When the user query is received, the appropriate copy is retrieved from disk and sent to the user. This first approach, often called *dynamic adaptation*, suffers from a very high CPU overhead for transcoding from one quality to another [23]. Therefore online transcoding is difficult in a multi-user environment. Our experiments (Fig 1) run on a 2.4GHz Pentium 4 CPU confirm this claim: a MPEG1 video is transcoded at a speed of only 15 to 60 frames per second. This corresponds to 60 - 240% of the entire CPU power if the frame rate for the video

<sup>1</sup>Another possible approach, which is not considered in this paper, is to use multilayered coding standards (MCS), which encode media into a base layer and multiple enhanced layers. QoS is provided by choosing what and how much data in the enhanced layers are delivered. An attractive alternative to the methods of transcoding and caching, MCS, however, will not totally replace them because it is only adapted by some coding formats such as MPEG4. Plus, the number of qualities MCS provides is still limited.

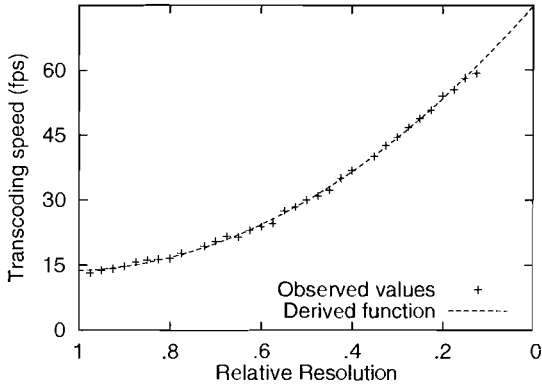


Figure 1: Time for transcoding a  $640 \times 480$  MPEG1 video to various (lower) resolutions.

is 25 frames per second (fps). We can see that CPU power is the bottleneck if we depend on online transcoding. As a result, many *transcode proxy* servers or video gateways [1] with massive computing power have to be deployed. The second approach, often called *static adaptation*, attempts to solve the problem of high CPU cost of transcoding by storing precoded multi-quality copies of the original media on disk. By this, the heavy demand on CPU power at runtime is alleviated. We trade disk space for runtime CPU cycles, which is a cost-effective trade-off since disks are relatively cheap.

Existing static adaptation systems are designed under one or both of the following assumptions: 1) *user requests concentrate on a small number of quality profiles*<sup>2</sup>; and 2) *there is always enough storage space*. However, these are not true for real-world multimedia databases. First of all, users vary widely in their quality needs and resource availability [23]. This leads to a large number of quality-specific copies of the same media content that need to be stored on disk. Secondly, although cheap, storage space is not free. This is especially true for commercial media databases that must provide high reliability of disk resources (which may be leased from vendors such as Akamai). Therefore, although storage is cheap, the storage requirements should not grow unboundedly. An analysis in Section 4 shows the disk space needed to accommodate all possible qualities could be intolerably high. Therefore, the choice of which quality copies to store becomes important and is the focus of this paper.

We view the selection of media copies for storage as a data replication problem (Fig. 2). Traditional data replication focuses on placement of copies of data in various nodes in a distributed environment [25]. Our quality-aware repli-

<sup>2</sup>As a result of 1), many media service providers offer quality options based on the client-side devices' processing capabilities. For example, CNN.com used to provide video streaming service in three different predefined qualities: one for dial-up users, one for DSL users, and for T1 users. However, this solution places strong limitations to the freedom of quality selection. Furthermore, with the development of mobile technologies, there are a large number of devices such as smart phones and PDAs, each of which has different rendering and communication capabilities. Thus, even by adapting this strategy, the number of possible quality-specific copies of the same media is large and keeps increasing with the emergence of new devices.

cation of multimedia deals with data placement in a metric space of quality values (termed as *quality space*). In the traditional replication scheme, data are replicated as exact or segmental copies of the original while the replicas in our problem are multi-quality copies generated via transcoding. In this paper, we present strategies to choose quality of replicas under two different user requirements: *Hard-Quality* and *Soft-Quality*. Under the hard-quality model, users must receive the exact quality requested. If such a quality is not already stored on disk, it must be generated by transcoding from an available quality. If the resources necessary for this transcoding are not available (e.g. due to too many requests) then the request is rejected. In a soft-quality model, users are willing to negotiate the quality that they receive and may be willing to accept a quality that is close to the original request. Naturally, there is a loss in utility for the user when he has to accept a different quality, depending upon the difference in quality. In either model, a request can be rejected if the system is overloaded (at the CPU, disk, or network).

Important performance metrics for these systems include: *reject rate* of requests, *user satisfaction*, and *resource consumption* [4, 17]. Our data replication algorithms are designed to achieve the lowest rejection rate, or highest user satisfaction under fixed resource (CPU, bandwidth, and storage) capacities. We hope our work will provide useful guidelines to system designers in building cost-effective and user-friendly multimedia databases.

The remainder of this paper is organized as follows: we first compare our work with others in Section 2; we then introduce the system model in Section 3; Section 4 discusses storage use of the replication process; we present our replica selection algorithms in Sections 5, 6, and 7; Section 8 is dedicated to experimental results; we conclude the paper by Section 9.

## 2. RELATED WORK AND OUR CONTRIBUTIONS

This work is motivated by the efforts to build quality-aware media systems [10, 23, 26, 30]. In our previous work [30], quality-aware query processing is studied in the context of multimedia databases. In that paper, we extend the query generation/optimization module of a multimedia DBMS to handle quality in queries as a core DBMS functionality. Two other related works in multimedia databases discuss quality specification [3] and quality model [32]. None of the above deals with replication of copies with different qualities. The closest work in quality-aware data replication is by Steinmetz *et al.* [27]. They focus more on availability and consistency of non-media data.

The traditional data caching/replication problem has been studied extensively in the context of web [28, 29], distributed databases [22, 25], and multimedia systems [19, 33]. The web caching and replication problem aims at higher availability of data and load balancing at the web servers. Similar goals are set for data replication in multimedia systems. What differs from web caching is that disk space and I/O bandwidth are the major concerns in multimedia systems. A number of algorithms are proposed to achieve high acceptance rate and resource utilization by balancing the use of different resources [6, 9, 33]. Unlike web and multimedia data, database contents are accessed by both read and write operations.

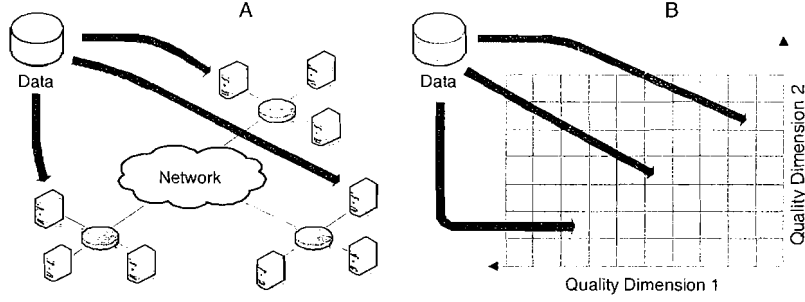


Figure 2: Traditional (A) and quality-aware (B) data replication.

This leads to high requirements on data consistency, which often conflict with data availability. Due to resource constraints, data consistency can sometimes only be enforced loosely.

Dynamic replication of data is another important issue. Access frequency to individual data items are likely to change in most environments. The goal is to make the replication strategy quickly and accurately adapt to changes and achieve optimal long-term performance. Wolfson *et al.* [34] introduced an algorithm that changes the location of replicas in response to changes of read-write patterns of data items. The interactions between query optimization and data cache utilization in distributed databases are discussed in [16]. They found that to take advantages of cached data, it is sometimes necessary to process individual queries using ‘sub-optimal’ plans in order to reach higher system performance. In [18] and [4], video replication/de-replication is triggered as a result of changes of request rates.

Quality support in media delivery in response to heterogeneous client features and environmental conditions has attracted a lot of attention [23, 24]. The problem of quality selection under storage constraints, however, has not been well addressed due to the oversimplified assumption of unlimited storage. A work that is close in spirit to ours is presented in [20] where quality selection is performed with the goal of minimizing expected transcoding costs between qualities. In another related paper [11], the problem of optimal materialized view selection is studied. Both [20] and [11] address different data selection problems from ours. Furthermore, neither considers quality selection in response to dynamic changes of query pattern. Another feature of our problem is that storage is shared by multiple physical objects, each of which has its own quality space. How to distribute storage among these quality spaces brings an extra dimension of difficulty. In comparison to the preceding papers, we make the following contributions in our study:

1. We analytically and experimentally show that the storage cost of static adaptation is so high that typically only a small number of replicas in the quality space can be accommodated in disks;
2. In a hard-quality system where users are assumed to be strict on quality requirements, we develop a (near-optimal) replica selection algorithm that minimizes request reject rate based on probabilistic analysis;
3. We formulate the replica selection under a soft-quality model as a facility location problem with the goal of

maximizing user satisfaction. We propose a fast greedy algorithm with performance comparable to commercial optimizers. An improvement to the greedy algorithm is also discussed; and

4. We extend the algorithms developed in 2) and 3) to handle dynamic changes of query pattern. Our solutions are fast and achieve the same level of optimality as the original algorithms. To the best of our knowledge, this is the first work to study the dynamic version of a combinatorial optimization problem.

A preliminary version of this paper appears in [31].

### 3. SYSTEM MODEL, NOTATIONS, AND ASSUMPTIONS

We assume that the database consists of a collection of servers that host the media content and service user queries. For now, we consider a centralized, single server scenario. The case of multiple, distributed servers is discussed in Section 6.5.3. We list in Table 1 the notations that will be used, throughout this paper.

Table 1: Notations and definitions.

Symbol	Definition
<b>System parameters</b>	
$B$	Server bandwidth
$S$	Server storage space for storing media
$C$	Server CPU power
$V$	Number of media objects in the system
$P$	Request reject rate
$M_i$	Total number of quality points for media $i$
$M$	Total number of quality points for all media
$f$	Total query rate, $f = \sum_{k=1}^M f_k$
<b>Quality-specific parameters</b>	
$f_k$	Query rate, number of requests per unit time
$\mu_k$	Service rate, requests served per unit time
$\lambda_k$	Request intensity, $\lambda_k = f_k \mu_k^{-1}$
$c_k$	CPU cycles per unit time for transcoding into this quality point from original quality
$b_k$	Bandwidth needed for streaming
$s_k$	Storage space needed if a replica is placed
$P_k$	Reject rate of requests to quality $k$

In our model, a server is characterized by the total amounts of the following resources available: bandwidth ( $B$ ), storage

space ( $S$ ), and CPU cycles ( $C$ ). Among them, *bandwidth* can be viewed as the minimum of the network bandwidth and the I/O bandwidth. In modern media servers, network bandwidth is most likely to be the bottleneck.

The system contains  $V$  media objects. User requests identify (either directly or via a query) an object to be retrieved as well as the desired quality requirements on  $m$  quality dimensions ( $\vec{q} = \{q_1, q_2, \dots, q_m\}$ , termed as *quality vector*). Each quality vector can thus be modeled as a point (hereafter called *quality point*) in a  $m$ -dimensional space. Generally, the domain of a quality parameter contains a finite number of values. For example, the spatial resolution of a video is an integer number of pixels within the range of  $192 \times 144$  (low-quality MPEG1) to  $1920 \times 1080$  (HDTV). Furthermore, the (horizontal) resolution can only be multiplies of 16 as the latter is the finest granularity most transcoders can handle. The total number of quality points for a specific media object  $i$  is  $M_i = \prod_{j=1}^m |Q_{ij}|$  where  $Q_{ij}$  is the set of possible values in dimension  $j$  for object  $i$  and  $Q_{ij}$  need not to be identical for all media objects. Note that every quality point is a candidate replica to be stored on disk.

Consider each possible quality,  $k$ , stored in the database. We use the following parameters to model this object:  $f_k, \mu_k, c_k, s_k, b_k$ .  $f_k$  represents the query rate for this version of the video. We assume that the query arrival is a Poisson process with this arrival rate. The query processing duration is assumed to follow an arbitrary distribution with expectation  $1/\mu_k$ . Note  $1/\mu_k$  may not be the same as the standard playback time of the media as the users may use VCR functionalities (e.g. stop, fast forward/backward) during media playback. The last three parameters ( $c_k, s_k, b_k$ ) correspond to the usage of resources. They can be precisely estimated from empirical functions derived by regression (see Section 4). Note  $c_k$  is fixed as the transcoding cost only depends on the target quality.

Under the hard-quality model, the server performs the following steps upon receiving a request:

- 1) attempts to retrieve from disk a replica that matches the quality vector  $\vec{q}$  attached to the request;
- 2) if the corresponding replica does not exist, transcodes a copy from a high-quality replica (by consuming  $c_k$  units of CPU) at runtime;
- 3) rejects the request if not enough CPU is available to perform 2).

If either 1) or 2) is performed, the retrieved/transcoded media data is transmitted to the client via the network (using  $b_k$  units of bandwidth). The request is also rejected if sufficient bandwidth is unavailable. We ignore the CPU costs of non-transcoding operations as they are trivial compared to transcoding costs and do not change with the specified  $\vec{q}$ . In the above model, requests are either admitted or rejected without waiting in a queue. The steps performed in soft-quality systems are slightly different from the above. We will discuss those in Section 6.

### 3.1 Assumptions

In this paper, we assume that replicas are readily available. In practice, all replicas can be precoded and archived on tertiary storage and copied into disk when a replication decision is made. We also assume that rejected queries are

**Table 2: Total relative storage in a 3D space.**

$n$	5	10	15	20	25
Storage	20.23	117.7	354.8	755.9	1496.5

not re-issued by the users. For analysis in Section 5, we make the following assumptions:

**Assumption 1.** We assume that CPU is a heavily overloaded resource as a result of online transcoding requests, i.e.,  $\sum_{k=1}^M \lambda_k c_k = m_c C$  where  $m_c \gg 1$ . On the other hand, the load put on system bandwidth is not as heavy as that on CPU, i.e.,  $\sum_{k=1}^M \lambda_k b_k = m_b B$  and  $m_b \ll m_c$ . We call  $m_b$  and  $m_c$  the *load coefficients* of these resources. Note the load on system bandwidth can be *critical*, i.e.,  $m_b = 1$  or *light*, i.e.,  $m_b < 1$ .

**Assumption 2.** We further assume that  $\frac{\min\{c_k\}}{B} > \frac{\max\{b_k\}}{C}$ , which means that the ratio of CPU cost to total CPU power is higher than that of bandwidth cost to total bandwidth for requests to all qualities.

The above two assumptions are reasonable due to our discussion in Section 1 about CPU being the bottleneck in our system model.

## 4. STORAGE REQUIREMENTS FOR QUALITY-AWARE REPLICATION

As mentioned in Section 1, it is often assumed in previous works that sufficient storage is available for static adaptation. Now we explore this assumption. Since a user can request any of the possible qualities, an ideal solution is to store most, if not all, of these replicas on disk such that only minimal load is put on the CPU for transcoding. We show that the storage cost for such a solution is simply too high.

We use digital video as an example throughout this paper. According to [24], the bitrate of a video replica with a single reduced quality parameter (e.g., resolution) is expressed as:

$$F = F_0(1 - R^\beta) \quad (1)$$

where  $F_0$  is the bitrate of the original video,  $R$  is the percentage of quality change ( $0 \leq R \leq 1$ ) from the original media, and  $\beta$  is a constant derived from experiments ( $0.5 < \beta < 1$ ). Suppose we replicate a media into  $n$  copies with a series of quality changes  $R_i$  ( $i = 1, 2, \dots, n$ ) that cover the domain of  $R$  evenly (i.e.  $R_i = i/n$ ). The sum of the bitrate of all copies is given by:

$$\sum_{i=0}^n F_0(1 - R_i^\beta) = F_0 \left[ n - \sum_{i=0}^n \left( \frac{i}{n} \right)^\beta \right] \quad (2)$$

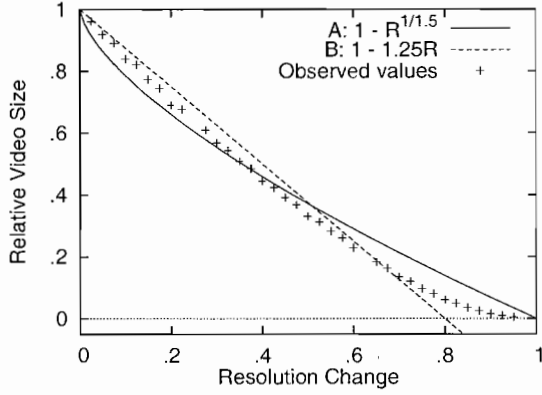
$$\approx F_0 \left[ n - \int_0^n \left( \frac{i}{n} \right)^\beta di \right] \quad (3)$$

$$= F_0 \left( n - \frac{n}{\beta + 1} \right) \quad (4)$$

$$= F_0 \frac{n\beta}{\beta + 1} = F_0 O(n). \quad (5)$$

The corresponding storage requirement can be easily calculated as  $T F_0 \frac{n}{\beta + 1}$  where  $T$  is the playback time of the media. Note that the above only considers one quality dimension. In [24], Equation (1) is also extended to three dimensions (spatial resolution, temporal resolution, and SNR):

$$F = \alpha F_0 (1 - R_a^\beta) (1 - R_b^\gamma) (1 - R_c^\delta) \quad (6)$$



**Figure 3: Change of video bandwidth with resolution degradation.**

where  $R_a$ ,  $R_b$ , and  $R_c$  are quality change in the three dimensions, respectively. The constants of their transcoder(s) are:  $\alpha = 1.12$ ,  $\beta = 2/3$ ,  $\gamma = 0.588$ , and  $\theta = 1.0$ . Using the same technique of approximation by integration as used in Equation (3), we can easily see the sum of all storage needed for all  $n^3$  replicas is  $TF_0O(n^3)$ . To be more general, the relative storage (to original size) needed for static adaptation is on the order of total number of quality points. The latter can be represented as  $O(n^d)$  where  $d$  is the number of and  $n$  is the replication density along quality dimensions. The conclusion is true as long as storage decreases polynomially with degradation of quality. Some of the storage costs generated using Eq. (6) are listed in Table 2. For example, when  $n = 10$ , the extra storage needed for all replicas is 117.7 times that of the original media size. No media service can afford to acquire hundreds of times of more storage for the extra feature of static adaptation. Needless to say, we could have even more quality dimensions in practice.

We have also experimentally verified the storage requirements for replication. We use the open-source video processing tool named *transcode*<sup>3</sup> in these experiments. Figure 3 shows the relative video size when spatial resolution decreases by various percentages. The discrete points are the resulting video sizes and curve A represents Eq. (1). In this graph, the areas under the curves can be viewed as the total relative storage use. We also plot a straight line B with function  $1 - 1.25R$  to show the theoretical storage usage based on Eq. (5). The area of the triangle formed by X, Y axes and line B is  $\frac{2}{5}$ , which is the same as that given by Eq. (5) since  $\beta = \frac{2}{3}$ . The fact that the areas under these three curves are very close to each other corroborates our analysis in this section.

## 5. HARD-QUALITY SYSTEMS

In this section, we discuss data replication strategies in hard quality systems where users have rigid quality requirements on service. This means a user is not willing to negotiate when the quality she specifies cannot be satisfied. As mentioned in Section 1, the main idea of static adaptation is to replicate original media into multiple quality copies such that the demand on CPU decreases. In Section 4 we have

shown that it is impractical to store all possible quality combinations. Therefore, the problem becomes how to choose quality points for replication given finite storage space  $C$  such that system performance is maximized. Since the non-availability of a requested quality results in the rejection of the request, we use the *reject probability*,  $P$ , as the metric for performance evaluation. Let the output of the replica selection algorithm be a vector  $(r_1, r_2, \dots, r_M)$  with 0/1 elements ( $r_k = 1$  if replica  $k$  is to be stored in disk). Formally, the replica placement problem is to

$$\begin{aligned} & \text{minimize } P, \\ & \text{subject to } \sum_{k=1}^M r_k s_k \leq S \end{aligned}$$

where the  $f_k$ ,  $\mu_k$ ,  $s_k$ , and  $c_k$  values for each replica are given.

To approach the above problem, it is critical to derive the relationship between  $P$  and the replica-specific values. First of all, the reject probability of all quality points is a weighted average of those of individual points:

$$P = \sum_{k=1}^M \frac{f_k}{\sum_{k=1}^M f_k} P_k = \frac{1}{f} \sum_{k=1}^M f_k P_k \quad (7)$$

where  $P_k$  is the reject probability of replica  $k$ . Suppose, by applying our replication algorithm, the  $M$  quality points are divided into two disjoint sets: a set  $\mathcal{R}$  containing replicated points and a set  $\mathcal{R}'$  with non-replicated points. Following Equation (7), we have

$$P = \frac{1}{f} (f_{\mathcal{R}} P_{\mathcal{R}} + f_{\mathcal{R}'} P_{\mathcal{R}'}) \quad (8)$$

where  $f_{\mathcal{R}} = \sum_{i \in \mathcal{R}} f_i$  is the total request rate in set  $\mathcal{R}$ ,  $P_{\mathcal{R}} = \frac{1}{f_{\mathcal{R}}} \sum_{k \in \mathcal{R}} f_k P_k$  is the reject probability of all requests from  $\mathcal{R}$  and  $f_{\mathcal{R}'}$ ,  $P_{\mathcal{R}'}$  are the counterparts of  $f_{\mathcal{R}}$ ,  $P_{\mathcal{R}}$  in set  $\mathcal{R}'$ .

In our model, the admission of a request is determined by the runtime availability of two resources: bandwidth and CPU. If either is insufficient to serve the request, the request is rejected. So the reject probability for a set of objects, say, those in set  $\mathcal{R}$ , can be expressed as

$$P_{\mathcal{R}} = P_{\mathcal{R}}^{(b)} + P_{\mathcal{R}}^{(c)} - P_{\mathcal{R}}^{(bc)} \quad (9)$$

where  $P_{\mathcal{R}}^{(b)}$ ,  $P_{\mathcal{R}}^{(c)}$ , and  $P_{\mathcal{R}}^{(bc)}$  are probabilities of the following events happening to requests from set  $\mathcal{R}$ : *rejected by bandwidth*, *rejected by CPU*, and *rejected by both CPU and bandwidth*. Note we cannot say  $P_{\mathcal{R}}^{(bc)} = P_{\mathcal{R}}^{(b)} \cdot P_{\mathcal{R}}^{(c)}$  as the first two events could be dependent on each other. Similarly, we have the following for set  $\mathcal{R}'$ :

$$P_{\mathcal{R}'} = P_{\mathcal{R}'}^{(b)} + P_{\mathcal{R}'}^{(c)} - P_{\mathcal{R}'}^{(bc)} \quad (10)$$

where  $P_{\mathcal{R}'}^{(b)}$ ,  $P_{\mathcal{R}'}^{(c)}$ , and  $P_{\mathcal{R}'}^{(bc)}$  are defined according to requests from set  $\mathcal{R}'$ .

As no rejection by CPU will occur when a replica is placed in disk (Section 3.1), we have  $P_{\mathcal{R}}^{(c)} = 0$ , which leads to  $P_{\mathcal{R}}^{(bc)} = 0$  and thus  $P_{\mathcal{R}} = P_{\mathcal{R}}^{(b)}$ . Plugging this and Eq. (10) into Eq. (8), we have

$$P = \frac{1}{f} (f_{\mathcal{R}} P_{\mathcal{R}}^{(b)} + f_{\mathcal{R}'} P_{\mathcal{R}'}^{(b)}) \quad (11)$$

$$= \frac{1}{f} [f_{\mathcal{R}} P_{\mathcal{R}}^{(b)} + f_{\mathcal{R}'} (P_{\mathcal{R}'}^{(b)} + P_{\mathcal{R}'}^{(c)} - P_{\mathcal{R}'}^{(bc)})] \quad (12)$$

We now establish the following proposition and theorem that will help analyze the above expression. Although both

<sup>3</sup><http://www.theorie.physik.uni-goettingen.de/~streich/transcode/>

Proposition 1 and Theorem 1 seem intuitively obvious, their proofs are non-trivial extensions of well-established results in queueing theory [5]. The basic idea is to map the hard-quality system to an Erlang loss model: we can view the bandwidth (CPU) as a resource pool with  $B$  ( $C$ ) channels, the replica-specific requests are modeled as Poisson streams with arrival rate  $f_k$ , and service rate  $\mu_k$ , and no waiting queue exists (i.e., lossy system). What complicates our analysis is that each request class requires a different number of channels (i.e.,  $c_k$ ,  $b_k$ ) whereas exactly one channel is used for one request in a regular Erlang system (e.g., one line for each telephone call). It is reported that stationary distributions do exist for the reject probability in such systems. The main results of such studies can be found in Appendix A.

**PROPOSITION 1.** *Given two Erlang loss systems  $A$  and  $B$ , each has a number of traffic classes. A class  $i$  in group  $A$  is characterized by its arrival rate  $f_i$ , service time  $\mu_i$ , and number of channels needed for each connection  $a_i$ . Similarly, a class  $j$  in group  $B$  is by  $f_j$ ,  $\mu_j$ , and  $b_j$ . Let  $R_A$  and  $R_B$  be the total number of channels for system  $A$  and  $B$ , respectively. If 1)  $R_A > R_B$ ; 2)  $m_A \gg m_B$  where  $m_A$  and  $m_B$  are the load coefficients of systems  $A$  and  $B$  (i.e., Assumption 1 in Section 3.1); and 3)  $\min_{i \in A} \{a_i\}/R_A > \max_{j \in B} \{b_j\}/R_B$  (Assumption 2 in Section 3.1), then the reject probability in system  $A$  is greater than that of system  $B$ .*

**PROOF.** See Appendix B.  $\square$

**THEOREM 1.** *If the requested load on a resource in the hard-quality system is critical or light, we have  $P = O\left(\sqrt{\frac{2\beta}{N\pi f^2}}\right)$  where  $N$  is the scale of resource pool (i.e.,  $N = \Theta(B)$  for bandwidth and  $N = \Theta(C)$  for CPU), and  $\beta = \sum_{k=1}^M f_k \mu_k$ .*

**PROOF.** See Appendix C.  $\square$

The unreplicated set  $\mathcal{R}'$  and replicated set  $\mathcal{R}$  can be mapped to groups  $A$  and  $B$  in Proposition 1, respectively. Since storage is limited, replicating some qualities does not change the fact that CPU is still heavily loaded (i.e., condition 2) always holds). Thus, we have

$$P_{\mathcal{R}}^{(b)} < P_{\mathcal{R}'}^{(c)} \leq P_{\mathcal{R}'}$$

when both resources are overloaded. This also holds true when the bandwidth (i.e., group  $B$ ) has critical or light load because the reject probability is close to zero under such conditions (Theorem 1, it is easy to see that  $\beta < f^2$  and  $N$  is large in our problem). The second inequality in the above formula is given by the fact that  $P_{\mathcal{R}'}^{(b)} \geq P_{\mathcal{R}'}^{(bc)}$ .

Revisiting Eq. (11), as  $P_{\mathcal{R}}^{(b)} < P_{\mathcal{R}'}$ , no matter how we choose members of  $\mathcal{R}$  and  $\mathcal{R}'$ , a heuristic solution to the problem of minimizing  $P$  would be to maximize  $f_{\mathcal{R}}$  (or minimize  $f_{\mathcal{R}'}$  since  $f_{\mathcal{R}} + f_{\mathcal{R}'} = f$ ) subject to  $\sum_{k \in \mathcal{R}} s_k \leq S$ . In other words, the problem becomes the classic 0-1 Knapsack problem, which can in turn be solved by the following heuristic algorithm: we sort all possible qualities by their request rate per unit size ( $f_k/s_k$ ) and select those with the highest such values till the total storage is filled. The running time of this algorithm is  $O(M \log M)$ . In Appendix D, we show that the results obtained by such a heuristic are near-optimal when  $S \gg s_k$  for all  $k \in [1, M]$ , which is a safe assumption.

The above result is interesting in that it shows that  $f_k$  and  $s_k$  are the only factors we need to consider in quality selection even though the reject probability is also a function of  $\mu_k$ ,  $c_k$ , and  $b_k$  (Appendix A).

## 6. SOFT-QUALITY SYSTEMS

In hard-quality systems, replicas of the same media object are treated as independent entities: storing a replica with quality  $q_1$  does not help the requests to another with quality  $q_2$  as quality requirements are either strictly satisfied or the request is not served at all. However, users can generally tolerate some changes of quality [24] and the quality parameters specified by a user only represent the most desirable quality. If these parameters cannot be exactly matched by the server, they are willing to accept similar set of qualities. The process of settling down to a new set of quality parameters is called *renegotiation*. Of course, the deviation of the actual qualities a user gets from those he/she desires will have some impact on the user's viewing experience and the system should be penalized for that.

### 6.1 Utility Functions

We generally use *utility* to quantify user satisfaction on a service received [21]. For our purposes, *utility functions* can be used to map quality to utility and the penalty applied to the media service due to renegotiation is easily captured by *utility loss*. As utility directly reflects the level of satisfaction from users, it is the primary optimization goal in quality-critical applications [17]. We thus set the goal of our replica selection strategies to be maximizing utility. The server operations shown in Section 3 needs to be modified in soft-quality systems. For simplicity, we assume the 'renegotiation' process between client/server is instantaneously performed on the server side based on a simple rule: in case of a miss in step 1), the server always chooses a replica that yields the largest utility for the request to retrieve.

Figure 4 shows various types of utility functions for a single quality dimension. In general, utility functions are convex monotones (Fig 4A) due to the fact that users are always happy to get a high-quality service, even if the quality exceeds his/her needs [21]. This makes our replica selection a trivial problem: always keep the one with the highest quality. However, in a more realistic environment, the cost of the extra quality may be high as more resources have to be consumed (Section 1) on the client side. Thus excessively high quality negatively affects utility. Taking this into account, we propose a new group of utility functions in quality-aware media services: it achieves the maximal utility at a single point  $q^{desire}$  and monotonically decreases on both sides of  $q^{desire}$  along the quality dimension (Fig 4B). The pattern of utility decrease with change of quality can either be dramatic (a, b of Fig 4B) or uniform (c of Fig 4B). Note that the functions do not have to be symmetric on both sides of  $q^{desire}$ . The hard-quality model in Section 5 can be viewed as a special case: its utility function takes the value of 1 at  $q^{desire}$  and 0 otherwise. The functions mentioned above are for one single quality dimension only. The utility for a quality vector with multiple dimensions is generally given as a weighted sum of dimensional utility described above [17]. The weights of individual quality dimensions are also user-dependent.

### 6.2 Data Replication as an Optimization



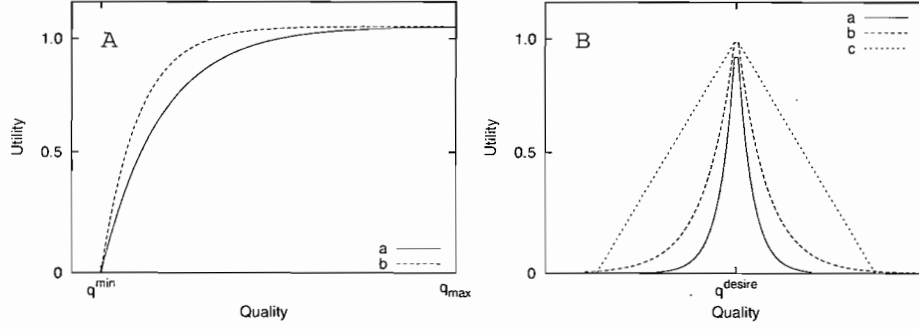


Figure 4: Different types of utility functions.

In this subsection we formally define the replica selection problem in soft-quality systems. Let us first study how to choose replicas for one media object  $i$ . We then extend our discussion to all  $V$  objects in Section 6.5.1.

The problem is to pick a set  $L$  of replicas that gives the largest total *utility* over time, which can be expressed as

$$U = \sum_{j \in J} f_j u(j, L)$$

where  $J$  is the set of all  $M_i$  points and  $u(j, L)$  is the largest utility with which a replica in  $L$  serves a request for quality  $j$ . Obviously,  $u(j, L)$  has maximum value when  $j \in L$ . We set  $u(j, L)$  to be a function of the distance between  $j$  and its nearest neighbor in  $L$  (Section 6.1). For example, when we put equal weights on both quality dimensions in a 2-D space and use linear functions such as  $c$  in Fig 4B,  $u(j, L)$  is actually the Manhattan distance between two points. Generally,  $u(j, L)$  is normalized into a value in  $[0, 1]$ . We weight the utility by the request rate  $f_j$  and the weighted utility is termed as *utility rate*. The constraint of forming set  $L$  is that the total storage of all members of  $L$  can not exceed  $S$ . We name our problem the *fixed-storage replica selection* (FSRS) problem and it can be formulated as the following integer programming:

$$\text{maximize} \quad \sum_{j \in J} \sum_{k \in J} f_j u(j, k) Y_{jk} \quad (13)$$

$$\text{subject to} \quad \sum_{k \in J} X_k s_k \leq S, \quad (14)$$

$$\sum_{k \in J} Y_{jk} = 1, \quad (15)$$

$$Y_{jk} \leq X_k, \quad (16)$$

$$Y_{jk} \in \{0, 1\}, \quad (17)$$

$$X_k \in \{0, 1\} \quad (18)$$

where  $u(j, k)$  is the utility value when a request to point  $k$  is served by a replica in  $j$ ,  $X_k$  is a binary variable representing whether  $k$  is replicated,  $Y_{jk}$  tells if  $j$  should be served by  $k$ . Equation (14) shows the storage constraint while Equations (15) and (16) mean all requests from  $k$  should be served by one and only one replica. Here  $f_j$ ,  $s_k$ , and  $S$  are inputs and  $X_k$  for all  $k \in J$  is the solution.

The FSRS problem looks similar to a group of well-studied optimizations known as the *uncapacitated facility location* (UFL) problems [7]. Yet it is different from any known UFL problems in that the storage constraint in FSRS is unique. A close match to FSRS is the so-called  $p$ -median problem with the same problem statements except Equation (14) becomes  $\sum X_k = p$ , meaning only  $p$  ( $p < |J|$ ) facilities are to be

built. As the  $p$ -median problem is NP-hard [12], we can thus conclude FSRS is also NP-hard.

**THEOREM 2.** *The FSRS problem is NP-hard.*

**PROOF.** The  $p$ -median problem is equivalent to finding the set of replicas that yields the smallest loss of utility rate in a quality space where  $s_k = \delta$  ( $\delta > 0$ ) for all  $k$  and  $S = p\delta$ . Thus the  $p$ -median problem is polynomial time reducible to the FSRS problem and this concludes the proof.  $\square$

### 6.3 The Greedy Algorithm.

As in the Knapsack problem, we can use a benefit/cost model to evaluate a replica  $k$ : the cost is obviously the storage  $s_k$ , the benefit would be the gain of utility rate of selecting  $k$ . What makes the problem more complicated is that the benefit is not fixed: it depends on the selection of other replicas. More specifically, the value of point  $k$  is the total utility rate of the set of points it serves and different selections of other replicas will affect the membership of this set of points. To bypass this difficulty, we propose an algorithm (we call it *Greedy*) that takes guesses on such benefits. The main idea is to aggressively select replicas one by one. The first replica is assigned to a point  $k$  that yields the largest  $\Delta U_k / s_k$  value as if only one replica is to be placed. We denote  $\Delta U_k / s_k$  as the *utility density* of replica  $k$  where  $\Delta U_k$  is the marginal utility rate gained by replicating  $k$ . The following replicas are determined in the same way with the knowledge of replicas that are already selected. The utility density value represents our guess of the benefit-to-cost ratio in replicating  $k$ . It should be noted that this is different from choosing the replicas simply in descending order of precomputed  $\Delta U_k / s_k$  because the  $\Delta U_k$  changes depending upon which replicas have already been selected and thus must be recomputed after each selection.

Fig. 5 shows the pseudo-code of the *Greedy* algorithm. GREEDY calls ADD-REPLICA continuously with a queue *list* holding the replicas selected so far. The algorithm terminates when no more replicas can be added due to storage constraints. The subroutine ADD-REPLICA is the core of this algorithm: it selects a new replica given those chosen in *list*. It does so by trying all  $M_i$  points in the quality space (line 2) to look for the one that yields the largest utility rate. Subroutine MAXUTILITY gives the utility from  $j$  to its nearest replica in *list* +  $k$ , which can be done in constant time if we store the previous nearest replica for all  $j$ . The two loops both have to run  $M_i$  iterations therefore the time complexity for *Greedy* is  $O(IM_i^2)$  for one media  $i$ . Here  $I$

**Algorithm GREEDY**  
**Inputs:**  $f_k, s_k (\forall k)$ , and  $S$   
**Output:** a set of selected replicas,  $P$

```

1  $s' \leftarrow S, P \leftarrow \emptyset, k \leftarrow 0$ 
2 while  $k \neq \text{NULL}$  do
3    $k \leftarrow \text{ADDREPLICA}(s', P)$ 
4    $s' \leftarrow s' - s_k$ 
5   append  $k$  to  $P$ 
6 return  $P$ 

```

**ADDREPLICA** ( $s, P$ )

```

1  $i \leftarrow \text{NULL}, V_{\max} \leftarrow 0$ 
2 for each quality point  $k$  do
3   if  $k \notin P$  and  $s_k \leq s$ 
4      $U \leftarrow 0$ 
5     for each quality point  $j$ 
6        $U \leftarrow U + \text{MAXUTIL}(j, k, P)$ 
7     if  $U/s_k > V_{\max}$ 
8        $V_{\max} \leftarrow U/s_k$ 
9      $i \leftarrow k$ 
10 return  $i$ 

```

Figure 5: The *Greedy* algorithm.

is the number of replicas eventually placed in *list*. In the worst case when all points are selected,  $I = M_i$ . In our storage constrained system,  $I$  should be asymptotically smaller than  $M_i$ .

**Effects of the type of utility functions.** It is easy to see that the shape of the utility functions affect the final results of replica selection. Recall that we evaluate a replica  $k$  by its  $\sum f_{ju}(j, k)/s_k$  value where  $j$  are the points  $k$  serves (line 7 in ADD-REPLICA). If the utility drops very fast, a replica can only collect utility from points that are extremely close to it therefore the *Greedy* algorithm favors those with high query rates in their close neighborhood. On the other hand, if utility drops very slowly, we may overestimate the utility rate of a point at early stages of *Greedy*. As a result, the first few replicas chosen by *Greedy* tend to be those with small  $s_k$  values since the utility rate of all candidates have little difference at that moment. In Section 6.4, we propose a solution to remedy this problem of *Greedy*.

The utility curves we have discussed so far are all monotonically decreasing functions of distance (between two points). However, our FSR algorithm does not depend on any special features (e.g. monotonicity, convexity) of the utility functions. In fact, *Greedy* works for arbitrary types of utility functions as long as the utility value between two points is not affected by the replica selection process.

## 6.4 The *Iterative Greedy* Algorithm

*Iterative Greedy* algorithm attempts to improve the performance of *Greedy*. We notice that at each step of *Greedy*, some local optimization is achieved: the  $(K + 1)$ -th replica chosen is the best given the first  $K$  replicas. The problem is: we do not know if the first  $K$  replicas are good choices. However, we believe the  $(K + 1)$ -th replica added is more

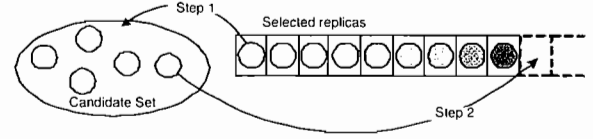


Figure 6: Replica replacement in *Iterative Greedy*.

**Algorithm ITERATIVEGREEDY**  
**Inputs:** selected replicas *slist*, number of iterations  $I$   
**Output:** a modified list of replicas *newlist*

```

1 copy slist to newlist
2  $U_{\max} \leftarrow 0, \text{storage} \leftarrow \text{available storage}$ 
3 for  $i \leftarrow 0$  to  $I$ 
4   do  $k \leftarrow \text{head of } slist$ 
5      $\text{storage} \leftarrow \text{storage} + s_k$ 
6      $l \leftarrow \text{ADD-REPLICA}(\text{storage}, slist)$ 
7     append  $l$  to slist
8     update storage,  $U \leftarrow \text{total utility of } slist$ 
9     if  $U_{\max} < U$ 
10        $U_{\max} \leftarrow U$ 
11     copy slist to newlist
12 return newlist

```

Figure 7: The *Iterative greedy* algorithm.

‘reliable’ than its predecessors because more global information (existence of other selected replicas) is leveraged in its selection. In this sense, the first replica is the most ‘unreliable’ one: it is chosen taking no such global information into account. Based on this conjecture, we develop the *Iterative Greedy* algorithm that iteratively improves the ‘correctness’ of the replicas chosen. Specifically, we repeatedly get rid of the most ‘unreliable’ selected replica and choose a new one, as illustrated in Fig 6. Note that the one that is eliminated is also a candidate of the selection process.

The operations in *Iterative Greedy* are shown in Fig 7. All replicas selected by *Greedy* are stored in a FIFO queue *slist*. In each iteration, we dequeue *slist* and find one replica (line 6) among the remaining replicas. The newly identified replica is then added to the tail of *slist*. The same subroutine, ADD-REPLICA, is used to find new replicas. Note that ADD-REPLICA may sometimes return no new replica if the one removed from *slist* leaves too little storage. We keep dequeuing *slist* and running ADD-REPLICA until  $I$  iterations are finished. We record the set of replicas with the largest utility rate as the final output. As ADD-REPLICA runs in  $O(M_i^2)$  time, *Iterative Greedy* has time complexity of  $O(IM_i^2)$ , which is the same as that of *Greedy*. The only problem here is how to set the number of iterations  $I$ . Since the primary goal of *Iterative Greedy* is to reconsider the selection of the first few ‘unreliable’ replicas, we can set  $I$  to be smaller than the total number of replicas selected by *Greedy*.

## 6.5 Other issues

### 6.5.1 Handling Multiple Media Objects

With very few modifications, both *Greedy* and *Iterative Greedy* algorithms can handle multiple media objects. The idea is to view the collection of  $V$  physical media as replicas of one virtual data object. The different content in the physical media can be modeled as a new quality dimension called

*content*. A special feature of *content* is its lack of adaptability. i.e., any replica of the movie *Matrix* cannot be used to serve a request to the movie *Shrek*. Assume all physical media have a quality space with  $\hat{M}$  points. the FSRS problem with  $V$  media can be solved by simply running the *Greedy* algorithm for the virtual media with  $V\hat{M}$  points. Knowing that there is no utility gain between two replicas with different *content*, we only need to run the second loop (line 5) in ADD-REPLICA for those with the same content. Thus, the time complexity of GREEDY becomes  $O(IV\hat{M}^2)$ .

Note some quality parameters for physical media objects also lack adaptability. Video format is a good example. Without degradation of bitrate, replication along these quality dimensions requires even more storage than adaptable dimensions. However, it can reduce the time complexity of GREEDY the same way as the *content* dimension does.

### 6.5.2 Relaxing/tightening constraints

Another point is that we set a constraint of replicating at least one copy for the video in Eq. (15). In the multi-object scenario, we can further relax or tighten this constraint. To relax it, we allow no replica being selected for a video, modifying Eq. (15) to  $\sum_{k \in J} Y_{jk} \leq 1$ . This requires no changes to our algorithms. On the other hand, the system administrator could also enforce the selection of certain replicas (e.g. the original video). Again, our FSRS algorithms can easily handle this: we just start running ADD-REPLICA with a list of all replicas that *must* stay in storage. However, if we stick to the original constraint but do not specify which replica to store for each video, the problem becomes trickier as our algorithms may assign no replica to videos with low query rates. The solution is to start by selecting the smallest replica for all videos and run *Greedy*. This guarantees one replica for each video but the effects of the constraint are minimized. Unless specified otherwise, the following extensions are based on a multi-video environment with the relaxed constraint.

### 6.5.3 Distributed Data Replication

In Section 5 and 6 we discuss the strategies of quality-aware data replication in a single server. Now we extend the solutions to a distributed system with multiple servers. Let us first investigate how the problem is changed when we consider multiple servers in a hard-quality system. Here we assume user requests can be served by any one of  $n > 1$  servers<sup>4</sup>. As we can see, the analysis we show in Section 4 still holds true and we can use Eq. (11) to guide our replica selection: the strategy is, again, to maximize  $f_R$ . When we obtain a set of replicas with the largest possible  $f_A$ , how to assign these replicas to  $n$  servers becomes a problem. We can immediately relate this to a load balancing problem with the goal of achieving uniform reject probability in all servers. A more detailed justification can be found in Appendix E.

For example, we can utilize a load balancing approach based on the idea of *resource pricing* proposed in [2]. In this approach, we set a *price* for the resource on which load is placed in each server. The price is set to reflect the supply-demand relationship of the resource. In our problem, the

price for bandwidth can be set to

$$\psi_{bandwidth} = n^{B'/B} \quad (19)$$

where  $B'$  is the load put on bandwidth so far. The replica placement is accomplished by putting replicas one by one into a server with the lowest cost. Note that in our problem we need to balance both storage and bandwidth. Therefore, the cost of placing replica  $k$  in a server is:

$$Cost = s_k \psi_{storage} + \lambda_k b_k \psi_{bandwidth}. \quad (20)$$

Resource prices are updated upon placement of each replica according to Eq. (19). The advantages of this algorithm are: server capacities do not have to be identical and it is proved to be  $O(\log n)$ -competitive [2].

The same strategy can be deployed to balance load under the soft quality system model even though reject probability is not the primary optimization goal.

## 7. DYNAMIC DATA REPLICATION

In previous sections we considered the situation of *static* data replication, in which access rates of all qualities do not change over time. The importance of studying static replication can be justified by two observations: 1. Access patterns to many media systems, especially video-on-demand systems, remain the same within a period of at least 24 hours [19]; 2. Conclusions drawn from static replication studies form the basis of dynamic replication research [18]. In this section, we discuss quality-aware data replication in an environment where access patterns change. There are two main requirements to a dynamic replication scheme: quick response to changes and optimality of results. Our goal is to design real-time algorithms that match static replication algorithms in terms of result optimality.

### 7.0.4 hard-quality systems

Our replication strategy for hard-quality systems is easily adaptable to dynamic situations: the replication decision is made by sorting replicas by their  $\eta_k = f_k/s_k$  values. When the query rate of a replica changes, we just reinsert the replica into the sorted list and make decisions based on its current position in the list. The algorithm is displayed in Fig 8. Recall from Section 5, all replicas belong to either the replicated set  $A$  or the non-replicated set  $B$ . In HARD-DYNAREP, we set a bound  $\bar{\eta}$  such that for any replica  $k$ , we have  $\eta_k > \bar{\eta} \Leftrightarrow k \in A$  and  $\eta_k < \bar{\eta} \Leftrightarrow k \in B$ . HARD-DYNAREP is called when we detect a change of access rate for a replica  $r$ . Replication decision is made based on comparison between the new  $\eta_r$  and the bound  $\bar{\eta}$ . The time complexity of this algorithm is  $O(\log M)$ .

### 7.0.5 Soft-quality systems

Dynamic replication in soft-quality systems is a very challenging task. The difficulty comes from the fact that the access rate change of a single point could have cascading effects on the choice of many (if not all) replicas. We may have to rerun the static algorithms (e.g. *Greedy*) in response to such changes but these algorithms are too slow to make online decisions. Fortunately, the *Greedy* and *Iterative Greedy* algorithms we developed have properties that we can exploit in building efficient, accurate dynamic replication algorithms. In this section, we assume that runtime variations of access pattern only exist at the media object level. In other words,

<sup>4</sup>If each server only handles requests from its local region, the problem is not interesting as we only need to perform single-server replication at each server.

```

Algorithm HARDDYNAREP
Inputs: sorted list  $L$  of all replicas,  $\bar{\eta}$ , and a replica  $r$ 
1  reinsert  $\eta_r$  into  $L$ 
2  Case 1.  $f_r$  increases
3    Case 1.1.  $r$  was replicated, do nothing
4    Case 1.2.  $r$  was not replicated
5      if  $\eta_r > \bar{\eta}$ 
6        do reset bound  $\bar{\eta}$ 
7         $\forall k$ , if  $k \in A$  and  $\eta_k < \bar{\eta}$ , dereplicate  $k$ 
8        replicate  $r$ 
9  Case 2.  $f_r$  decreases, operations are opposite to Case 1.

```

Figure 8: hard-quality dynamic replication algorithm.

the relative popularities of different quality points for the same media object do not change. Although this assumption is reasonable in many systems [19, 33], we understand a solution for more general situation is meaningful and we leave it as future work.

## 7.1 Replication roadmap

Let us first investigate how ADD-REPLICA, being the core of both *Greedy* and *Iterative Greedy*, selects replicas. The history of total utility rate gained and storage spent on each selected replica can be represented as a series of points in a 2D graph. We call the lines that connect these points in the order of their being selected a *Replication Roadmap* (RR). Fig. 9 shows two examples of RRs plotted with the same scale. We can see that any RR is convex. The reason for this is: the slope of the line connecting any two consecutive points (e.g.  $r_1$  and  $r_2$  in Fig. 9) in a RR represents the ratio of  $\Delta U_{r_2}$  to  $s_{r_2}$ . As ADD-REPLICA always chooses a replica with the largest  $\Delta U/s$  value, the slopes of the lines along the RR are thus non-increasing.

We can also draw RRs for individual media objects. It is not hard to see that single-media RRs are also convex. In dynamic replication, replicas need to be re-selected with respect to the new query rate of a media object. Suppose the query rate  $f_i$  of a medium  $i$  increases by a factor  $\delta$  ( $\delta > 0$ ). This makes the slopes of all pieces in  $i$ 's RR increase by  $\delta$ . What happens now is that we may consider assigning extra storage to  $i$  as it reaches a position to use storage more profitably than before. As storage is limited, the extra chunk should come from another medium whose slope in the last piece of RR is small. Take Fig 9 as an example. Suppose we have fully extended RRs: all future replicas are precomputed (empty dots in Fig 9) and we call the last real replica the *frontier* of the RR. It buys us more utility to advance A's frontier (take storage) and move backwards on B's RR (give up storage). The beauty of this scheme is that: we never need to pick up points far into or over the frontier to make storage exchanges. The convexity of RRs tells us that the frontier is always the most efficient point to acquire/release storage. Based on this idea, we have the an online algorithm named SOFTDYNAREP for dynamic replication (Fig. 10).

## 7.2 The SOFTDYNAREP algorithm

The algorithm consists of two phases: the *Preprocess Phase* and *Online Phase*. In the Preprocess phase, we need to extend each RR formed by *Greedy* or *Iterative Greedy* by

```

Algorithm SOFTDYNAREP
Preprocess Phase
1  run GREEDY or Iterative Greedy
2  for all  $V$  media objects
3    store the post-frontier segment of the RR in  $flist$ 
4    store the pre-frontier segment of the RR in  $blist$ 
5    extend RR to its full length
Online Phase
6  Case 1.  $f_i$  increases
7    recalculate slopes of stored segments for  $i$ 
8    update  $blist$  and  $flist$  (reinsertion)
9    do STORAGEEXCHANGE
10 Case 2.  $f_i$  decreases, symmetric to Case 1

STORAGEEXCHANGE
1   $storage \leftarrow$  available storage
2   $k \leftarrow 0, j \leftarrow V - 1$ 
3  while  $k \leq 0$ 
4    do  $r_0 \leftarrow flist[k]$ 
5     $victims \leftarrow \emptyset$ 
6    while  $storage < \text{size of replica } r_0$ 
7      do  $r_1 \leftarrow blist[j]$ 
8      if  $r_0$  and  $r_1$  belong to the same video
9         $j \leftarrow j - 1$ 
10     continue
11     if utility density of  $r_1 > \text{utility density of } r_0$ 
12        $k \leftarrow k + 1$ 
13       rollback  $blist$  to its status on line 6
14     break
15     else append  $r_1$  to  $victims$ 
16     update and sort  $blist$ 
17   if  $storage \geq \text{size of replica } r_0$ 
18     EXCHANGE ( $r_0, victims$ )
19     update and sort both  $flist$  and  $blist$ 

```

Figure 10: Soft-quality dynamic replication algorithm.

adding all  $M_i$  replicas<sup>5</sup>. For all RRs, we put the immediate predecessor of the *frontier* in a list called *blist* and the immediate successor in a list called *flist*. Both lists are sorted by the slopes of the segments stored. The *Preprocess phase* runs at  $O(VM^3)$  time and it only needs to be executed once.

The *Online Phase* is triggered once we detect a change in query rate to an object  $i$ . The idea is to iteratively take storage from the end of *blist* until a new equilibrium is reached. The running time of this phase is  $O(I \log V)$  where  $I$  is the number of storage exchanges (line 9). In the worst case where most of  $i$ 's replicas are to be stored, we have  $I = O(M_i)$ . The case of query rate decrease is just handled in an opposite way to what we have discussed above. In EXCHANGESTORAGE, there are two loops: in the outer loop (line 3), we choose the replica ( $r_0$ ) on the head of *flist* and try to find a list (*victims*) of replicas on the tail of *blist* from where storage can be taken via the inner loop (line 6). The list *victims* has to be formed as the size of  $r_0$  can be larger than that of one single victim replica  $r_1$ . The subroutine EXCHANGE basically dereplicates those in *victims* and replicate  $r_0$ . The inner loop terminates when enough storage is found for  $r_0$  or we reach a replica whose utility density is greater than that of  $r_0$  (line 11). The latter case also terminates the outer loop (as  $k > 0$ ).

Suppose the access rate to a video increases by a factor of  $\delta$  and we rerun the *Greedy* algorithm, we shall see that the replicas are chosen following the same RR as before. This is

<sup>5</sup>In practice, we do not have to extend a RR to its full length if we can bound the possible changes of query rates.

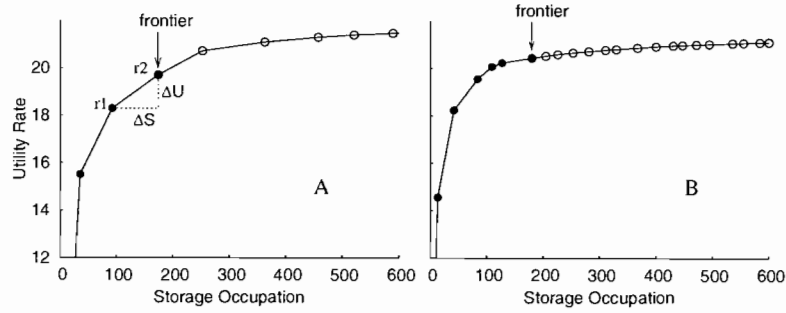


Figure 9: Replication roadmaps.

because the access rates to all quality points for this media increase proportionally as a result of the stable intra-media access pattern.

### 7.3 Optimality of SOFTDYNAREP.

In this section, we show that the online phase of SOFTDYNAREP achieves (almost) the same quality in the selected replicas as that by rerunning GREEDY at runtime.

From discussions in Section 7.1, we know that the global RR changes as the query rates of individual replicas change and GREEDY (implicitly) rebuilds the whole global RR. Essentially, *Greedy* selects those replicas with the largest utility density on the global RR, similar to our solution to the hard-quality problem (i.e., a 0-1 Knapsack). Let us consider a modified version of *Greedy* named M-GREEDY, which has a subtle difference from the original GREEDY algorithm. In M-GREEDY, we replicate items along the global RR till we encounter the first replica  $k'$  that cannot be accommodated by the available storage (equivalent to  $l$  in Appendix D)<sup>6</sup>. We immediately see that GREEDY is smarter than M-GREEDY: it will try to fill the available storage with replicas with lower utility density than  $k'$ . Thus, the replicas selected by M-GREEDY is a prefix of the global RR (from the beginning to the one prior to  $k'$ ) while those selected by GREEDY is not a consecutive chunk in the RR. Due to the same reasons discussed in Appendix D, the total utility rate achieved by M-GREEDY is only trivially smaller than that of GREEDY. To accomplish our claim that SOFTDYNAREP is as good as GREEDY, we have the following lemma.

**LEMMA 1.** *With the same replica-specific inputs and change of query rate of a specific video, if a replica is selected by M-GREEDY, it is also selected by SOFTDYNAREP.*

**PROOF.** Let us first study the change of the global RR before and after the query rate change. In Fig. 11, the global RR is represented as an array of replicas sorted by descending order of utility density. We know that GREEDY selects replicas from the left to the right till no storage is available. We draw a line called *boundary* between those that are replicated and those that are not. We consider the case of query rate increase of an object  $v$ . As a result of query rate increase, some replicas of  $v$  (represented as shaded boxes in Fig. 11) will move toward the left in the array of replicas and a new boundary will be formed. However, the relative order of all replicas of  $v$  does not change. Therefore, there are

<sup>6</sup>We use M-GREEDY as a conceptual variance to GREEDY, its implementation is irrelevant to our discussions.

two types of selected replicas by the M-GREEDY algorithm after the change: 1. those that were not selected before the change, and 2. those that were selected before the change. We prove SOFTDYNAREP selects the corresponding replicas in both cases:

**Case 1.** Without loss of generality, we consider a replica  $k$  of  $v$  that moves across the boundary in M-GREEDY. The selection of  $k$  can be achieved by one of two means: 1. the storage left before the change is greater than  $s_k$ ; 2. storage is taken from replicas with utility density smaller than that of  $k$ . It is easy to see that  $k$  will be the head of *flist* in SOFTDYNAREP. In the former case, we directly go to line 17 of STORAGEEXCHANGE (Fig. 10) and replicate  $k$ . For the second situation, a list of replicas are chosen to give up their storage to  $k$  (loop in line 6). As long as there is enough storage from those with smaller utility density,  $k$  will be replicated.

**Case 2.** The replicas considered in this case can be divided into two categories:

**Case 2.1.** Replicas whose utility density is greater than that of  $k$  (e.g., those in region  $S_0$  in Fig 11). These replicas are not affected by SOFTDYNAREP as we never sacrifice such replicas for  $k$  (line 11 of STORAGEEXCHANGE, Fig. 10).

**Case 2.2.** Replicas whose utility density is smaller than that of  $k$  (e.g., those in region  $S_1$  in Fig 11). These replicas are part of region  $S$  before the query rate changes. One feature of M-GREEDY is that all replicas chosen form a consecutive chunk in the list. To accommodate  $k$ ,  $S$  is simply cut into two consecutive regions  $S_1$  and  $S_2$ . In SOFTDYNAREP, the same list *victims* is also a consecutive chunk as it is formed by always choosing the replica with the smallest utility density, starting (backwards) from the end of  $S$ . Furthermore, it ends as long as enough storage is found thus everything in  $S_1$  will not be included in *victims*.

The case of multiple replicas crossing the boundary and decrease of query rate would not complicate the above argument.  $\square$

## 8. EXPERIMENTS

We study the behavior of various algorithms described in previous sections by extensive simulations. We use traces of 270 MPEG1/2 videos extracted from a real video database as experimental data<sup>7</sup>. For all replicas, we set their  $\mu_k$  to be their standard playback time. The videos are then transcoded into replicas of different spatial resolution and

<sup>7</sup><http://www.cs.purdue.edu/vdbms>

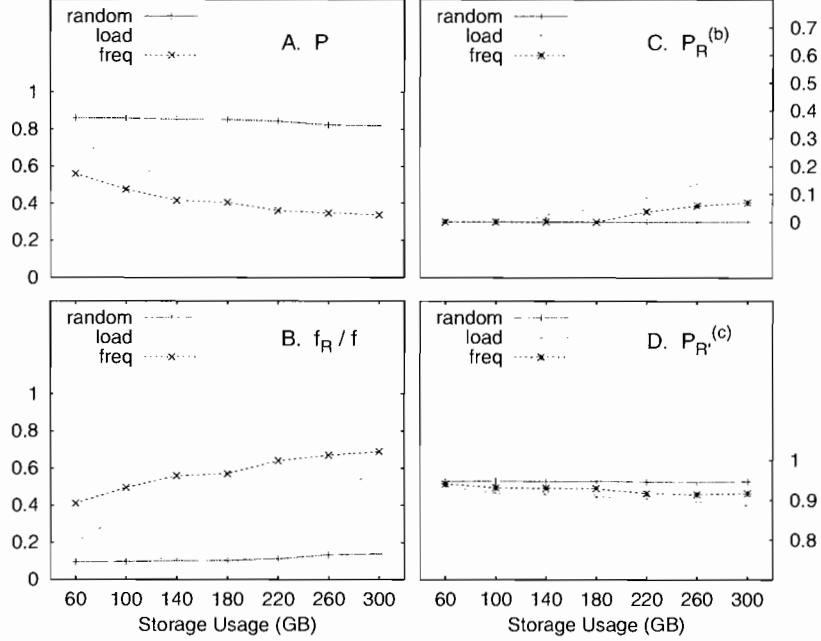


Figure 12: Performance of various replica selection algorithms in the hard-quality model.

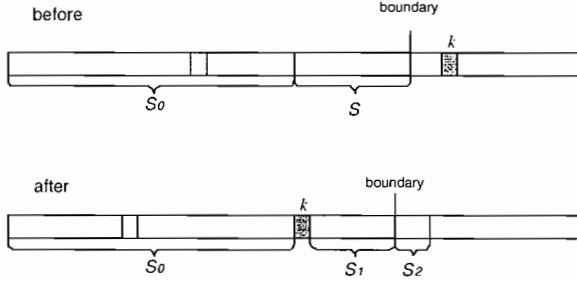


Figure 11: Replica selection upon query rate change.

frame rates using *transcode* (Section 4) to generate the  $b_k$ ,  $c_k$ , and  $s_k$  values for all replicas. We test various access patterns (e.g., uniform, Zipf, 20-80, 10-90) in our simulations. The simulated video servers possess network bandwidth of 90Mbps (dual T3 lines), four UltraSparc 1.2MHz CPUs, and variable storage capacity (60 to 300G) for data replication. All the above parameters are set to be close to those in a real-world server<sup>8</sup> and we simulate a cluster of 10 such servers. We perform our simulations on a Sun Workstation with a UltraSparc 1.2MHz CPU and 2 gigabytes of min memory running Solaris 8.

## 8.1 Results for Hard-Quality Model

In this experiment, we compare our replica selection algorithm (Section 5) to various heuristics under the hard quality system model. The metric is the reject frequency measured

<sup>8</sup>The total storage is relatively small because we only have 270 raw videos in the simulated system and the quality of most of the videos are not very high. In real systems, storage is more abundant but we may also have much more raw media with higher quality.

as the ratio of the total number of rejected requests to total requests. The quality space is a 2-D space (resolution and frame rate) with 15 to 20 values on each dimension (differs by each video object). Requests (with  $f = 7200/\text{hour}$ ) are distributed in a Zipf pattern to all  $M$  replicas.

In Fig. 12, we show the performance of three quality selection methods: 1) our solution that chooses quality points by their  $f_k/s_k$  values ('*freq*'); 2) an algorithm that randomly chooses replicas one by one till all storage is filled ('*random*'); and 3) an algorithm ('*load*') that places the largest possible load into set  $\mathcal{R}$  by choosing qualities with the largest CPU load to storage ratio ( $\frac{\lambda_k c_k}{s_k}$ ). The results confirm our analysis in Section 5 as our solution (*freq*) always gets the lowest reject probability (Fig 12A). As expected, the total request rate in set  $\mathcal{R}$  achieved by *freq* is always the highest (Fig 12B). In fact, the  $f_{\mathcal{R}}$  value achieved in our results are very close to those given as the upper bound of the optimal such values in all cases<sup>9</sup>.

From Fig. 12C and Fig. 12D we can see that the rejection frequency on bandwidth is significantly smaller than that on CPU. In these experiments, the recorded load coefficient of bandwidth ( $m_b$ ) range from 0.24 to 1.26. On the other hand, the load coefficients of CPU rendered by the same jobs range from 16 to 36. This explains why the observed reject frequency on CPU (Fig. 12D) is always high ( $> 0.88$ ). For algorithms *freq* and *load*, as storage increases,  $P$  and  $P_{\mathcal{R}}^{(c)}$  decrease while  $P_{\mathcal{R}}^{(b)}$  and  $f_{\mathcal{R}}$  increase. Note when excessive storage is used the decrease of  $P$  slows down as congestion on bandwidth becomes more significant. The performance of *random* is not affected by total storage.

## 8.2 Results for Soft-Quality Model

<sup>9</sup>The upper bound is described in Appendix D and is not plotted here.

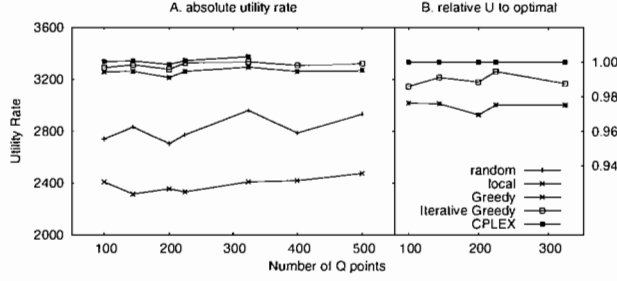


Figure 13: Optimality of replica selection algorithms.

In this section we present experimental results under the soft-quality model. We first evaluate the performance of GREEDY and ITERATIVEGREEDY algorithms in terms of optimality (Fig. 13) and running time (Fig. 14). In this experiment, we set  $f$  to 3600 requests/hour thus the utility rate is bounded by 3600/hr. We compare our algorithms with three others: 1) the CPLEX mathematical programming package<sup>10</sup>. CPLEX is a widely-used software for solving various optimization problems and is well-known for its efficiency. We tune CPLEX such that the utility rate of results obtained are within a 0.01% gap to that of the optimal solution; 2) the same *random* algorithm as the one described in Section 8.1; 3) a *local* algorithm that places replicas in the hottest areas in the quality space<sup>11</sup>. We run the experiments for a total of 30 media objects<sup>12</sup>. Each data point represents the mean of four simulations.

From Fig 13A, it is clear that our algorithms always find solutions that are very close to the optimal. More details can be found in Fig 13B where the relative  $U$  values obtained by our algorithms to those by CPLEX are plotted. Utility rates of solutions found by GREEDY are only about 3% smaller than the optimal values. The ITERATIVE GREEDY cuts the gap by at least half in all cases: its solutions always achieve more than 99% of the optimal utility rate. The performance of both algorithms is not affected by the increase of number of quality points. Neither is it affected by access patterns: we tested different access patterns (e.g. Zipf, 20-80, 10-90, and uniform) and obtained similar results (data not plotted). The solutions given by *random* and *local* are far from optimal. Surprisingly, the *local* algorithm, which is similar to our solution under the hard-quality model (Section 5), performs even worse than the random algorithm. This reiterates that it is dangerous to consider only local or regional information in solving a combinatorial problem.

The running time of the above experiments are shown on a logarithmic scale in Fig 14. CPLEX is the slowest algorithm in all cases. This is what we expected as its target is always the optimal solutions. Actually, we could only run CPLEX for the five smaller cases due to its long run-

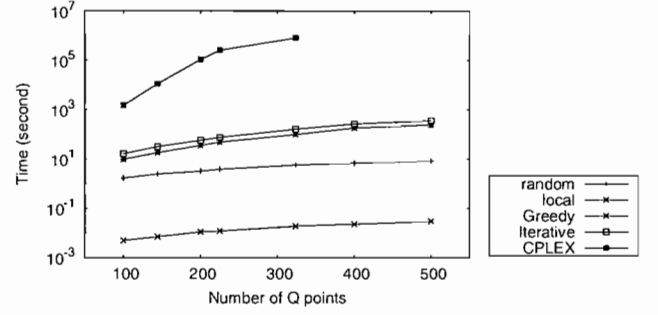


Figure 14: Running time of different replica selection algorithms.

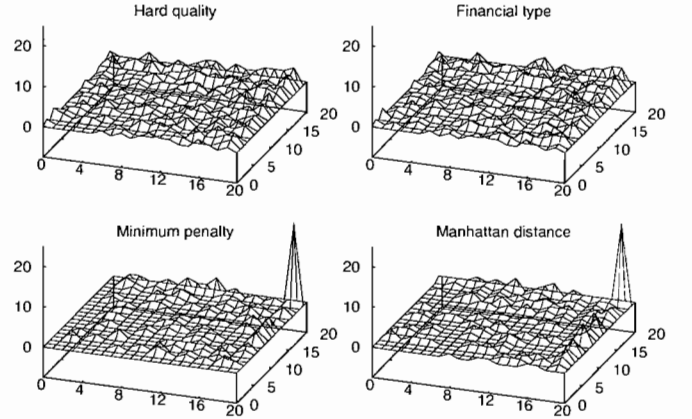


Figure 15: Frequency of replicas chosen by Greedy in a  $20 \times 20$  quality space.

ning time. Both *Greedy* and *Iterative Greedy* are 2-4 orders of magnitude faster than CPLEX. It takes them about 200 seconds to solve the selection of 30 videos in a quality space with 500 points. From Section 6 we know the running time increases linearly with the number of media objects  $V$ . Thus, it may take a few hours to select replicas in a real media system with thousands of media objects. Fortunately, we do not need to run these algorithms very often and the running time of our online algorithm is very small, as we will see in Section 8.3.

### 8.2.1 Effects of utility functions

We test our algorithms with four types of utility functions: *hard-quality*, *financial*, *Manhattan distance*, and *minimum penalty*. They are ordered by the speed of utility loss as a function of distance in the quality space. The type of utility function we used in the experiments presented in Fig. 13 and Fig. 14 is Manhattan distance. The details of these utility functions are as follows:

1. *hard-quality*, which only gains utility when the distance  $d$  between the requested quality and retrieved quality is zero, i.e.,  $U = \begin{cases} 1.0 & \text{if } d = 0 \\ 0 & \text{otherwise} \end{cases}$
2. *Financial type*, in which utility decreases exponentially with the distance  $d$  between two quality points, e.g.,  $U = 1.0 - e^{-d\tau}$  where  $\tau > 1$ ;
3. *Manhattan distance*, which is defined as the sum of

<sup>10</sup>version 8.0.1, <http://www.cplex.com>

<sup>11</sup>Specifically, we divide the whole space into squares of equal size. Each square is evaluated by the sum of the query rates of all quality points it contains. We place a replica in the center of those squares with the largest total query rates till all storage is filled.

<sup>12</sup>Here we choose a small number of media so that it is feasible to find the optimal solutions as a comparison to our solutions.



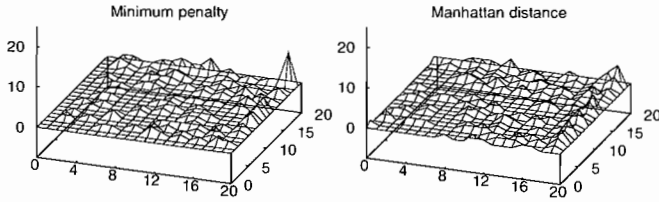


Figure 16: Frequency of replicas chosen by *Iterative Greedy* in a  $20 \times 20$  quality space.

dimensional differences between two points, e.g.,  $d = |x_1 - x_2| + |y_1 - y_2|$  for two points with coordinates  $(x_1, y_1)$  and  $(x_2, y_2)$ . In this type, utility decreases linearly with distance  $d$ .

4. Minimal penalty. In this type, utility remains 1.0 if the distance  $d$  is smaller than a threshold value and decrease linearly afterwards otherwise.

Fig. 15 shows the frequency of quality points chosen by *Greedy* in a  $20 \times 20$  space for a total of 30 videos. In Fig. 15, any point  $(x, y, z)$  shows that, out of the 30 media objects,  $z$  objects have replicas of quality  $(x, y)$  selected by the algorithm. Larger numbers on X, Y axes mean lower quality. We can see that utility functions significantly affect the choice of replicas. For *hard-quality* and *financial* whose utility drops very fast, the replicas are evenly distributed in the quality space. For the other two utility functions, *Greedy* selects more replicas with lower quality. A salient problem is that for over 20 videos, *Greedy* picks the lowest quality replica (19, 19). This confirms our discussion in Section 6.3: with overestimated utility rates, smaller replicas are always chosen first. The situation is improved by the *Iterative Greedy* algorithm. Fig 16 shows the distribution of replicas after running *Iterative Greedy* with the same set of inputs. The high peaks on points (19, 19) disappear and total utility rate increases by about 2%.

It should be noted that the solutions found by *Greedy* are almost optimal if we use *hard-quality* and *financial* types of utility functions. *Iterative Greedy* has no advantages under this situation. Our explanation to this is: by utilizing fast utility-dropping functions, we are making the FSRS problem a lot easier to solve. Recall (Section 6) that the major difficulty of solving FSRS comes from the combinatorial effects among replicas in collecting utility. However, the above utility functions tend to make replicas more isolated as they can only collect utility locally.

### 8.3 Dynamic Replication

We study replica selection in a multi-server environment under the hard-quality model. Experimental setup is the same as that described in Section 8.1 except the simulator contains 10 identical servers. We compare the performance of three strategies: load balancing by *resource pricing* (Section 6.5.3), load balancing by *Bandwidth-storage ratio* (BSR) [6], and random assignment of load. Fig 17A shows the results of load balancing using the metric of standard deviation normalized by the mean of loads. The *pricing* strategy has slightly better performance than BSR. The *random* method generates highly unbalanced load distribution. The effect of load balancing on reject rate is presented in Fig 17B: the *random* method performs the worst while *pricing* only

has marginal advantages over BSR. From this experiment we conclude that load balancing is needed. However, it is not clear which load balancing strategy is better and further investigation on this is beyond the scope of this paper.

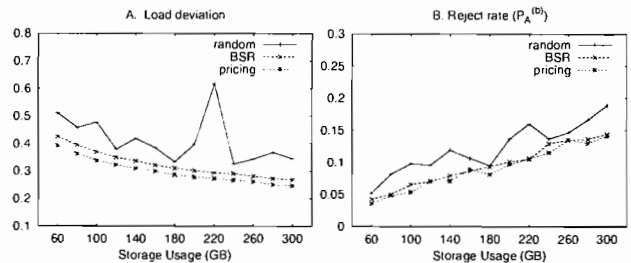


Figure 17: Performance of load balancing methods.

We also test our dynamic replication algorithm for the soft-quality model for its optimality and speed. We simulate a system for a period of time during which events of query rate changes of media objects are randomly generated. We allow the query rate of videos to increase up to 20 times and to decrease down to 1/10 of the original rate. We first compare the total utility rate of the selected replicas between the online phase of SOFTDYNAREP and *Greedy*. In all cases, the replicas selected match exactly with those found by the M-GREEDY algorithm discussed in Section 7.3 thus the utility rates are always the same between two solutions. As shown in Fig 18A, the replicas selected by SOFTDYNAREP have utility rates that are consistently within 99.5% of that by the original *Greedy* algorithm. In this experiment with 270 videos and a  $20 \times 20$  quality space, the running time of SOFTDYNAREP for each event is on the order of  $10^{-4}$  seconds while ADD-REPLICA needs to run about half a hour to solve the same problems. The main reason for SOFTDYNAREP's efficiency is the small number of storage exchanges. In Fig. 18B, we record such numbers for each execution of SOFTDYNAREP and very few of these readings exceeds 15. This shows that our algorithm is suitable for making real-time decisions.

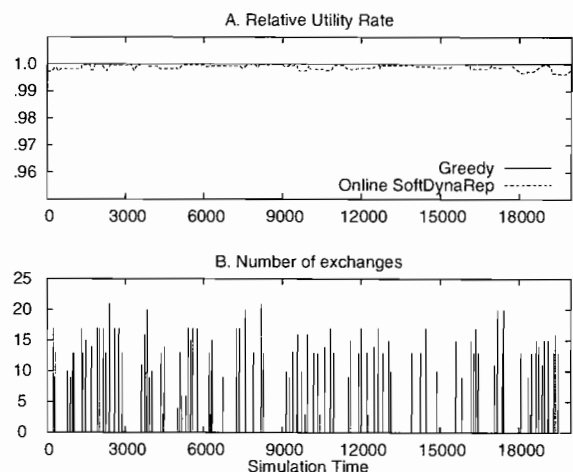


Figure 18: Performance of SOFTDYNAREP.

## 9. CONCLUSIONS



In this paper, we study the problem of selecting quality-specific replicas of media data. This problem is generally ignored in multimedia database research due to the oversimplified assumption that storage space is abundant. We demonstrate by analysis and experiments that this is not the case if the system is to adapt to user quality requirements with reasonable granularity. We provide solutions to the problem under two different system models. In the discussions on a hard-quality system model, we conclude the query rate and storage of individual replicas are the most critical factors that affect performance. We also propose a greedy algorithm to solve the replica selection problem on a soft-quality system model. Experiments show that the solutions found by our algorithm are within 3% of the optimal. An advanced version of this algorithm further reduces that to 1%. A derived online algorithm provides an elegant solution to an important subproblem of dynamic data replication.

## APPENDIX

### A. REJECT (BLOCKING) PROBABILITY IN A GENERALIZED ERLANG MODEL

We use our CPU as an example to elaborate this. Then the CPU requests from different replicas can be viewed as competitors for a shared resource pool with a finite capacity  $C$ . Kelly first studied the probability of rejection in such systems [14, 15]. The main idea is to analyze the occurrence of resource occupation states denoted as  $\vec{n} = (n_1, n_2, \dots, n_M)$  where  $n_k$  is the number of requests to replica  $k$  currently being serviced. According to [15], the reject probability of any replica  $k$  is

$$P_k = \frac{\sum_{\vec{n} \in S_k} \prod_{k=1}^M \frac{1}{n_k!} \lambda_k^{n_k}}{\sum_{\vec{n} \in \mathcal{S}} \prod_{k=1}^M \frac{1}{n_k!} \lambda_k^{n_k}} \quad (21)$$

where  $S_k = \{\vec{n} : C - c_k < \sum_{k=1}^M n_k c_k \leq C\}$  and  $\mathcal{S} = \{\vec{n} : \sum_{k=1}^M n_k c_k \leq C\}$  are two sets of states. The states in  $S_k$  are those at which a request to replica  $k$  will be rejected (as there are less than  $c_k$  units of resource available) while  $\mathcal{S}$  is the collection of all possible states.

Due to the discrete feature of the states, it is very difficult to discuss the characteristics of Eq. (21). Fortunately, Gazdzicki *et al.* [8] gives the following asymptotic approximation to Eq. (21).

**Case 1.** When the resource has light load, i.e.,  $\sum_{k=1}^M \lambda_k c_k < C$ , the class-specific reject probability is

$$P_k = e^{\tau d \epsilon - I(C)} \frac{d}{\sqrt{2\pi N \sigma}} \left( \frac{1 - e^{\tau c_k}}{1 - e^{\tau d}} \right) (1 + o(1)). \quad (22)$$

**Case 2.** When the resource has critical load, i.e.,  $\sum_{k=1}^M \lambda_k c_k = C$ ,  $P_k$  becomes

$$P_k = \sqrt{\frac{2}{\pi N}} \frac{c_k}{\sigma} (1 + o(1)). \quad (23)$$

**Case 3.** When the pool is heavily loaded, i.e.,  $\sum_{k=1}^M \lambda_k c_k > C$ , we get

$$P_k = (1 - e^{\tau c_k}) (1 + o(1)) \quad (24)$$

where  $N$  is the *scale* of resource pool (i.e.,  $N = \Theta(C)$ ),  $\tau$  is

the unique solution to the equation

$$\sum_{k=1}^M \frac{f_k}{\mu_k} c_k e^{\tau c_k} = C, \quad (25)$$

and other relevant quantities are defined as follows:

- i.  $d$  is the greatest common divisor of  $c_1, c_2, \dots, c_M$ ;
- ii.  $\epsilon = \frac{C}{d} - \lfloor \frac{C}{d} \rfloor$  where  $\lfloor a \rfloor$  denotes the largest integer such that  $\lfloor a \rfloor \leq a$ ;
- iii.  $I(C) = \tau C - \sum_{k=1}^M \lambda_k (e^{\tau c_k} - 1)$ ;
- iv.  $\sigma^2 = \sum_{k=1}^M \lambda_k c_k^2 e^{\tau c_k}$ .

### B. PROOF OF PROPOSITION 1

**PROOF.** The reject probability for group  $\mathcal{A}$  is  $P_A = \frac{1}{f_A} \sum_{i \in \mathcal{A}} f_i P_i$  where  $P_i$  is the reject probability for traffic class  $i$  and  $f_A = \sum_{i \in \mathcal{A}} f_i$ . Similarly, we have  $P_B = \frac{1}{f_B} \sum_{j \in \mathcal{B}} f_j P_j$ . For overloaded resources, we can use Eq. (24) to quantify the quality-specific reject probability. Therefore, we get  $P_i = 1 - e^{\tau_A a_i}$  and  $P_j = 1 - e^{\tau_B b_j}$  where  $\tau_A$  and  $\tau_B$  are constants that satisfy Eq. (25). Let  $s = \tau_A$  and  $u = \tau_B$  and we call  $s$  and  $u$  the *passage coefficients*<sup>13</sup> of groups  $\mathcal{A}$  and  $\mathcal{B}$ . To prove  $P_A > P_B$ , it is sufficient to show that  $\frac{1}{f_A} \sum_{i \in \mathcal{A}} f_i s^{a_i} < \frac{1}{f_B} \sum_{j \in \mathcal{B}} f_j u^{b_j}$ .

We first apply a proportional scaling to the classes in group  $\mathcal{A}$ , that is, we increase all  $a_i$  as well as the total resource units  $R_A$  by a factor of  $\omega$  ( $\omega > 1$ ) such that  $\omega R_A = R_B$ . According to [8], such “scaling” will not increase the class-specific reject probability, i.e.,  $\forall i, P'_i \leq P_i$  where  $P'_i = 1 - s'^{\omega a_i}$  is the reject probability of class  $i$  after it is scaled. Note that  $s$  is replaced by a new constant  $s'$ . With this transformation, this proof is concluded if we can show

$$\frac{1}{f_A} \sum_{i \in \mathcal{A}} f_i s'^{\omega a_i} < \frac{1}{f_B} \sum_{j \in \mathcal{B}} f_j u^{b_j} \quad (26)$$

Kelly [14] states that the value of the passage coefficient of a traffic group can be approximated by the inverse of its load coefficient. Thus, we get  $s' \approx \frac{1}{m_A}$  and  $u \approx \frac{1}{m_B}$ . Note that scaling does not change the load coefficient of group  $\mathcal{A}$ . For  $m_A \gg m_B$ , we have  $s' < u$  for sure. With the given condition  $\min_{i \in \mathcal{A}} \{a_i\}/R_A > \max_{j \in \mathcal{B}} \{b_j\}/R_B$ , we immediately have  $\min_{i \in \mathcal{A}} \{\omega a_i\} > \max_{j \in \mathcal{B}} \{b_j\}$ , which further leads to  $s'^{\omega a_i} < u^{b_j}, \forall i, j$ . Having this, formula (26) is trivially correct.  $\square$

### C. PROOF OF THEOREM 1

**PROOF.** Using the notations for the bandwidth resource  $(\lambda_k, b_k, B)$ , we first derive an upper bound of  $P$  under the critical load situation. Recall the asymptotic approximation to  $P_k$  for a critically-loaded resource in Eq. (23), we immediately have  $\tau = 0$  and  $e^{\tau b_k} = 1$ .

From Eq. (23), we get  $\lambda_k P_k^2 = \frac{2}{\pi N} \frac{\lambda_k b_k^2}{\sum_{k=1}^M \lambda_k b_k^2}$ , which leads to

$$\sum_{k=1}^M \lambda_k P_k^2 = \sum_{k=1}^M f_k \mu_k^{-1} P_k^2 = \frac{2}{\pi N}. \quad (27)$$

<sup>13</sup>These are basically the probabilities of one single unit of resource being free.

To get the upper bound for  $P = \frac{1}{f} \sum f_k P_k$ , we use the method of Lagrangian multipliers with the following optimization function

$$L = \sum f_k P_k - \phi \left( \sum f_k \mu_k^{-1} P_k^2 - \frac{2}{\pi N} \right)$$

where  $\phi$  is the Lagrangian multiplier. We discuss how  $P_k$  may affect the bound of  $P$  given all  $f_k$  and  $\mu_k$ . The condition for maximality is thus  $\frac{\partial L}{\partial P_k} = 0, \forall k$ . This is the same as

$$f_k - 2\phi \frac{f_k}{\mu_k} P_k = 0, \forall k.$$

Immediately, we get  $P_k = \frac{\mu_k}{2\phi}$  as the condition for achieving the upper bound. Plugging this into Eq. (27), we have  $2\phi = \sqrt{\frac{\pi N \sum f_k \mu_k}{2}}$ . Let  $\beta = \sum f_k \mu_k$ , we have  $P_k = \mu_k \sqrt{\frac{2}{\pi N \beta}}$  under the optimal situation. Therefore, the maximum value of  $P$  can be expressed as

$$P = \frac{1}{f} \sum f_k P_k = \frac{1}{f} \sum f_k \mu_k \sqrt{\frac{2}{\pi N \beta}} = \sqrt{\frac{2\beta}{\pi N f^2}}.$$

Now we consider the underload situation (i.e., Case 1 in Appendix A). Note that any such case can be transformed to a critical load case by adding a new class (i.e., class  $M+1$ ) of requests. Specifically, we let  $b_{M+1} = B$  and choose an arbitrary  $\lambda_{M+1}$  such that  $\sum_{k=1}^M \lambda_k b_k + \lambda_{M+1} b_{M+1} = B$ . We now show that the reject probability  $P$  always increases after the transformation.

According to [13],  $P_k$  can be obtained from the following relation:

$$\sum_{k=1}^{M+1} \lambda_k b_k q(j - b_k) = j q(j), \quad j = 0, 1, \dots, B \quad (28)$$

where  $q(j)$  is the stationary probability that exactly  $j$  units of resources are occupied and  $q(j) = 0$  for  $j < 0$ . It is easy to see that  $\sum_{j=0}^B q(j) = 1$ , and  $P_k$  is given by  $P_k = \sum_{i=0}^{b_k-1} q(B-i)$ .

Running Eq. (28) recursively with the unknown quantity  $q(0)$  as the base case, we have

$$q(0) + q(1) + q(2) + \dots + q(B) = q(0)(1 + \alpha_1 + \alpha_2 + \dots + \alpha_B) = 1 \quad (29)$$

where  $\alpha_j$  ( $1 \leq j \leq B$ ) is a constant determined by the recursions. By adding class  $M+1$ , reconsidering Eq. (28), Eq. (29) becomes

$$q(0) + q(1) + q(2) + \dots + q(B) = \quad (30)$$

$$q(0) [1 + \alpha_1 + \alpha_2 + \dots + (\alpha_B + \lambda_{M+1} b_{M+1})] = 1.$$

As a result, the value of  $P_k = \sum_{i=0}^{b_k-1} q(B-i)$  for any class  $k$  is larger in Eq. (30) than in Eq. (29). Therefore, quantity  $P$  in the underload case is smaller than the corresponding critical load case generated by the above transformation. In other words, it is also bounded by  $\sqrt{\frac{2\beta}{\pi N f^2}}$ .  $\square$

## D. OPTIMALITY OF A SIMPLE SOLUTION TO THE 0/1 KNAPSACK PROBLEM

In the 0-1 Knapsack problem, each candidate object has its own size and value. We define the ratio of value over size as the value density of an object (denoted as  $v$ ). We claim that if we put the objects with the largest value density

into the knapsack, the total value obtained are near-optimal when the size of the knapsack is far greater than the size of any individual object.

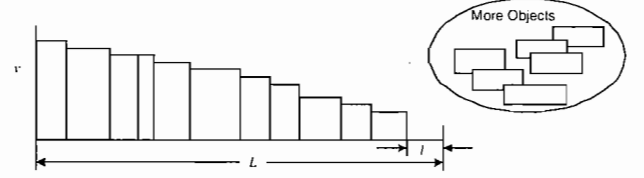


Figure 19: A knapsack filled with objects.

As illustrated in Fig. 19, the knapsack has size  $L$ , and  $Y$  axle represents value density. Each candidate object is represented as a rectangle and its size as the width, and value as the area of the rectangle, respectively. Our algorithm will fill the knapsack with objects with the largest value densities until no objects can be filled as a whole. It is easy to see that, if all  $L$  storage is filled, the solution is optimal as any other plan will decrease the total value achieved. If there is a unfilled region with size  $l$ , we can fill it with the object whose value density  $v'$  is the largest among the unselected objects. This generates an infeasible solution as we have to cut a piece (with size  $l$ ) from the last object. However, it gives an upper bound of the optimal total value:  $\hat{v} = \Omega + lv'$  where  $\Omega$  is the achieved total value (area of the shaded region in Fig. 19). Here  $lv'$  can be viewed as an upper bound of the difference between our solution and the optimal value. When  $L \gg l$ , we have  $lv' \ll \hat{v}$ .

## E. LOAD BALANCING IN DISTRIBUTED DATA REPLICATION

Suppose set  $\mathcal{R}$  contains certain number of replicas and they are to be placed on  $n$  servers. Denote the total query rate in server  $i$  as  $f_i$  and reject probability as  $P_i^{(b)}$ . Immediately, we have  $f_{\mathcal{R}} = \sum_{i=1}^n f_i$  and  $P_{\mathcal{R}}^{(b)} = \frac{1}{f_{\mathcal{R}}} \sum_{i=1}^n P_i^{(b)}$ . In distributed replication, we have the problem of *minimizing*  $P_{\mathcal{R}}^{(b)}$  given  $f_{\mathcal{R}}$ , which can be solved by using Lagrangian multipliers with the following solution:

$$\frac{\partial}{\partial f_i} \left[ \sum_{i=1}^n f_i P_i^{(b)} - \phi \left( \sum_{i=1}^n f_i - f_{\mathcal{R}} \right) \right] = 0, \quad \forall f_i \quad (31)$$

where  $\phi$  is a Lagrange multiplier. For any  $i$ , the LHS of Equation (31) is  $P_i^{(b)} - \phi$ . Thus, we get the following condition of minimality

$$P_1^{(b)} = P_2^{(b)} = \dots = P_n^{(b)} = \phi. \quad (32)$$

Theoretically, it is not clear how to achieve uniform  $P_i^{(b)}$  in our case. Little *et al.* [19] proved that, when all requests have the same bandwidth requirements ( $b_k$ ), the above condition is achieved when all servers have the same load ( $\sum \lambda_k b_k$ ) on bandwidth. In our system where requests have different bandwidth requirements, the above condition can only be approximated by evenly distributing the load (Theorem B is favorable to this approach).

## Acknowledgments

The authors would like to thank Prof. Xingquan Zhu and Mr. Shan Lei, for sharing with us their valuable insights in the early stages of this study. We are also grateful to Prof. Hong Wan, for her help with the CPLEX optimization software.

## F. REFERENCES

- [1] E. Amir, S. McCanne, and H. Zhaing. An Application Level Video Gateway. In *Proceedings of ACM Multimedia*, pages 255–265, 1995.
- [2] J. Aspens, Y. Azar, A. Fiat, S. Plotkin, and O. Waarts. On-Line Load Balancing with Applications to Machine Scheduling and Virtual Circuit Routing. In *Proceedings of ACM STOC*, pages 623–631, 1993.
- [3] Elisa Bertino, Ahmed Elmagarmid, and Mohand-Said Hacid. A Database Approach to Quality of Service Specification in Video Databases. *SIGMOD Record*, 32(1):35–40, 2003.
- [4] C.-F. Chou, L. Golubchik, and J. C. S. Lui. Striping Doesn't Scale: How to Achieve Scalability for Continuous Media Servers with Replication. In *Proceedings of IEEE ICDCS*, pages 64–71, April 2000.
- [5] Robert B. Cooper. *Introduction to Queueing Theory*. North Holland, New York, 1981.
- [6] A. Dan and D. Sitaram. An Online Video Placement Policy Based on Bandwidth to Space (BSR). In *Proceedings of ACM SIGMOD*, pages 376–385, 1995.
- [7] Z. Drezner and H. W. Hamacher. *Facility Location: Applications and Theory*. Springer, 2002.
- [8] P. Gazdzicki, I. Lambadaris, and R. Mazumdar. Blocking Probabilities for Large Multirate Erlang Loss Systems. *Advances in Applied Probability*, 25:997–1009, December 1993.
- [9] L. Golubchik, S. Khanna, S. Khuller, R. Thurimella, and A. Zhu. Approximation Algorithms for Data Placement on Parallel Disks. In *Proceedings of SODA*, pages 223–232, 1998.
- [10] A. Hafd and G. Bochmann. An Approach to Quality of Service Management in Distributed Multimedia Application: Design and Implementation. *Multimedia Tools and Applications*, 9(2):167–191, 1999.
- [11] V. Harinarayan, A. Rajaraman, and J. D. Ullman. Implementing Data Cubes Efficiently. In *Proceedings of SIGMOD*, pages 205–216, June 1996.
- [12] O. Kariv and S. L. Hakimi. An Algorithmic Approach to Network Location Problems. II: The  $p$ -Medians. *SIAM Journal of Applied Mathematics*, 37(3):539–560.
- [13] J. S. Kaufman. Blocking in a Shared Resource Environment. *IEEE Transactions on Communications*, 29(10):1474–1481, October 1981.
- [14] F. P. Kelly. Blocking Probabilities in Large Circuit-Switched Networks. *Advances in Applied Probability*, 18(2):473–505, June 1986.
- [15] F. P. Kelly. Loss Networks. *Annals of Applied Probability*, 1(3):319–378, August 1991.
- [16] D. Kossman, M. Franklin, and G. Drasch. Cache Investment: Integrating Query Optimization and Distributed Data Placement. *ACM Trans. of Database Systems*, 25(4):517–558, 2000.
- [17] C. Lee, J. Lehoczy, D. Siewiorek, R. Rajkumar, and J. Hansen. A Scalable Solution to the Multi-Resource QoS Problem. In *Proceedings of the IEEE Real-Time Systems Symposium*, December 1999.
- [18] P. W. K. Lie, J. C. S. Lui, and L. Golubchik. Threshold-Based Dynamic Replication in Large-Scale Video-on-Demand Systems. *Multimedia Tools and Applications*, 11:35–62, 2000.
- [19] T. D. C. Little and D. Venkatesh. Popularity-Based Assignment of Movies to Storage Devices in a Video-on-Demand System. *Springer/ACM Multimedia Systems*, 2(6):280–287, January 1995.
- [20] W. Y. Lum and F. C. M. Lau. On Balancing Between Transcoding Overhead and Spatial Consumption in Content Adaptation. In *Proceedings of MOBICOM*, pages 239–250, September 2002.
- [21] G. Menges. *Economic Decision Making: Basic Concepts and Models*, chapter 2, pages 21–48. Longman, 1973.
- [22] A. Milo and O. Wolfson. Placement of Replicated Items in Distributed Databases. In *Proceedings of EDBT*, pages 414–427, 1988.
- [23] R. Mohan, J. R. Smith, and C.-S. Li. Adapting Multimedia Internet Content for Universal Access. *IEEE Transactions on Multimedia*, 1(1):104–114.
- [24] S. Nepal and U. Srinivasan. DAVE: A System for Quality Driven Adaptive Video Delivery. In *Proceedings of Intl. Workshop of Multimedia Information Retrieval (MIR04)*, pages 224–230.
- [25] M. Nicola and M. Jarke. Performance Modeling of Distributed and Replicated Databases. *IEEE Trans. Knowledge and Data Engineering*, 12(4):645–672, July/August 2000.
- [26] B. Noble, M. Price, and M. Satyanarayanan. A Programming Interface for Application-aware Adaptation in Mobile Computing. In *Proc. 2nd USENIX Symposium on Mobile Computing*, April 1995.
- [27] Giwon On, Jens Schmitt, Michael Liepert, and Ralf Steinmetz. Replication with QoS support for a Distributed Multimedia System. In *Proceedings of the 27th EUROMICRO Conf., Warsaw, Poland*, September 2001.
- [28] L. Qiu, V. N. Padmanabhan, and G. M. Voelker. On the Placement of Web Server Replicas. In *Proceedings of IEEE INFOCOM*, pages 1587–1596, 2001.
- [29] M. Rabinovich. Issues in Web Content Replication. *Data Engineering Bulletin*, 21(4):21–29, 1998.
- [30] Y.-C. Tu, S. Prabhakar, A. Elmagarmid, and R. Sion. QuaSAQ: An Approach to Enabling End-to-End QoS for Multimedia Databases. In *Proceedings of EDBT*, pages 694–711, March 2004.
- [31] Yi-Cheng Tu, Jingfeng Yan, and Sunil Prabhakar. Quality-Aware Replication of Multimedia Data. In *Proceedings of Database and Expert Systems Applications (DEXA)*, pages 240–249, August 2005.
- [32] J. Walpole, C. Krasic, L. Liu, D. Maier, C. Pu, D. McNamee, and D. Steere. Quality of Service Semantics for Multimedia Database Systems. In *Proceedings of Data Semantics 8: Semantic Issues in Multimedia Systems*, volume 138, 1998.
- [33] Y. Wang, J. C. L. Liu, D. H. C. Du, and J. Hsieh.