

# Improving Multi-dimensional Query Processing with Data Migration in Distributed Cache Infrastructure

Youngmoon Eom<sup>†</sup>, Jinwoong Kim<sup>†</sup>, Deukyeon Hwang<sup>†</sup>, Jaewon Kwak<sup>†</sup>, Minho Shin<sup>‡</sup>, Beomseok Nam<sup>†</sup>

<sup>†</sup> Dept. of Computer Science Engineering  
Ulsan National Institute of Science and  
Technology (UNIST)  
Republic of Korea

<sup>‡</sup> Dept. of Computer Engineering  
Myongji University  
Republic of Korea

{youngmoon01,jwkim,deukyeon,sovioletta,bsnam}@unist.ac.kr, mhshin@mju.ac.kr

**Abstract**—In distributed query processing systems where caching infrastructure is distributed and scales with the number of servers, it is becoming more important to orchestrate and leverage a large number of cached objects in distributed caching systems seamlessly as the present trend is to build large scalable distributed systems by connecting small heterogeneous machines. With a large scale distributed caching system, a scheduling policy must consider both cache hit ratio and system load balance to optimize multiple queries. A scheduling policy that considers system load but not cache hit ratio often fails to reuse cached data by not assigning a query to the sever that has data objects the query needs. On the contrary, a scheduling policy that considers cache hit ratio but not system load balance may suffer from system load imbalance.

To maximize the overall system throughput and to reduce query response time, a multiple query scheduling policy must balance system load and also leverage cached objects. In this paper, we present a distributed query processing framework that exhibits high cache hit ratio while achieving good system load balance. In order to seamlessly manage our distributed scalable caching system, our framework performs autonomic cached data migrations to improve cache hit ratio. Our experiments show that our proposed query scheduling policy and data migration policy significantly improve system throughput by achieving high cache hit ratio while avoiding system load imbalance.

## I. INTRODUCTION

Multi-dimensional data processing has been the subject of extensive research in various fields including scientific applications, database systems, computer graphics, geographic information systems, etc. A large portion of real world datasets are multi-dimensional, and enormous amount of multi-dimensional datasets are generated and analyzed by modern computer systems.

Efficient processing of such large multi-dimensional datasets is a major challenge in many disciplines, and distributed and parallel query processing systems have been successfully used for scientific applications to solve such large

scale complicated data analysis problems because substantial performance gains can be obtained by exploiting data and computation parallelism.

When dealing with computationally intensive scientific applications, load balancing plays an important role in reducing query response time and improving system throughput via task parallelism. In addition to the load balancing, data intensive applications are known to benefit from reusing intermediate cached results if multiple queries can exploit sub-expression commonality [1].

In distributed query processing systems, the size of distributed caches scales with the number of distributed servers. Leveraging the large distributed caches plays an important role in improving overall system throughput but it is not an easy problem to reuse cached objects while maintaining load balance. Load-based scheduling policies do not consider cached objects in distributed caching system and get little benefits from reusing them. To take the advantages of large scale distributed caching infrastructure, more intelligent query scheduling policies are required. However, a cache-aware query scheduling policy does not always imply high system throughput because it may hurt system load balance. Suppose a scheduler knows a remote server has some popular cached objects, so it forwards most subsequent queries to the remote server just to improve cache hit ratio, but it will flood the server with a large number of queries while all the other servers are idle. Although in general cache hit ratio helps reducing query response time, if a server is flooded with too many waiting queries, the waiting time can be much higher than the benefit of cache hit.

In this paper, we study a distributed query scheduling policy - *DEMA* that makes query scheduling decisions based on *spatial locality* of multi-dimensional queries so that similar queries get clustered together in order to improve cache hit ratio and the number of queries in each cluster becomes balanced. Our proposed spatial clustering algorithm is different from the well known spatial clustering algorithms used in machine learning in a sense that the query scheduling decisions must be made at run time, i.e., the overhead of run-time scheduling algorithms should be minimal. *DEMA*

This research was supported by MKE/KEIT (No.10041608, Embedded System Software for New Memory based Smart Devices) and ICT R&D program of MSIP/IITP. (No. 1391104004, Development of Device Collaborative Giga-Level Smart Cloudlet Technology)

scheduling policies take into account the dynamic contents of distributed caching infrastructure in order to exploit sub-expression commonality of cached results and employ statistical prediction methods to balance load across distributed servers.

In order to manage the distributed caching infrastructure in a seamless fashion, our distributed query processing framework also employs an autonomic data migration policy. Because DEMA scheduling policy makes scheduling decision based on spatial locality, if the distribution of clustered cached objects changes and some cached objects do not belong to its cluster any more, our framework migrates the cached object to a remote server where the cached object can be reused. Our extensive experimental studies show our proposed query scheduling policy and its data migration policy significantly improves the query processing throughput and outperforms legacy load based scheduling policy by a large margin.

The rest of the paper is organized as follows: Section II describes other research works related to cache-aware query scheduling policies and query optimization in distributed systems. In section III, we describe overall architecture of our distributed query processing framework. In section IV, we discuss DEMA scheduling policy and how the scheduling policy achieves both load balance and high cache hit ratio. In section V, we propose a data migration policy for DEMA scheduling algorithm. In Section VI we evaluate our scheduling policy, and in section VII we conclude.

## II. RELATED WORK

Load-balancing problems have been extensively investigated in many different fields. Godfrey et. al [2] proposed an algorithm for load balancing in heterogeneous and dynamic peer-to-peer systems. Catalyurek et. al [3] investigated how to dynamically restore balance in parallel scientific computing applications where the computational structure of the applications change over time. Vydyanathan et. al [4] proposed scheduling algorithms that determine what tasks should be run concurrently and how many processors should be allocated to each task. Zhang et. al [5] and Wolf et al. [6] proposed scheduling policies that dynamically distribute incoming requests for clustered web servers. WRR (Weighted Round Robin) [7] is a commonly used, simple but enhanced load balancing scheduling policy which assigns a weight to each queue (server) according to the current status of its load, and serves each queue in proportion to the weights. However, none of these scheduling policies were designed to take into account a distributed cache infrastructure, but only consider the heterogeneity of user requests and the dynamic system load.

LARD (Locality-Aware Request Distribution) [8], [9] is a locality-aware scheduling policy designed to serve web server clusters, and considers the cache contents of back-end servers. The LARD scheduling policy causes identical user requests to be handled by the same server unless that server is heavily loaded. If a server is too heavily loaded, subsequent user requests will be serviced by another idle

server in order to improve load balance. The underlying idea is to improve overall system throughput by processing queries directly rather than waiting in a busy server for long time even if that server has a cached response. LARD shares the goal of improving both load balance and cache hit ratio with our scheduling policies, but LARD transfers workload only when a specific server is too heavily loaded while our scheduling policies actively predict future workload balance and take actions beforehand to achieve better load balancing.

In relational database systems and high performance scientific data processing middleware systems, exploiting similarity of concurrent queries has been studied extensively. That work has shown that heuristic approaches can help reuse previously computed results from cache and generate good scheduling plans, resulting in improved system throughput as well as reducing query response times [10], [1]. Zhang et al. [11] evaluated the benefits of reusing cached results in a distributed cache framework in a Grid computing environment. In that simulation study, it was shown that high cache hit rates do not always yield high system throughput due to load imbalance problems. We solve the problem with scheduling policies that consider both cache hit rates and load balancing.

In order to support data-intensive scientific applications, a large number of distributed query processing middleware systems have been developed including MOCHA [12], DataCutter [10], Polar\* [13], ADR [10], and Active Proxy-G [10]. Active Proxy-G is a component-based distributed query processing grid middleware that employs user-defined operators that application developers can implement. Active Proxy-G employs meta-data directory services that monitor performance of back-end application servers and a distributed cache indexes. Using the collected performance metrics and the distributed cache indexes, the front-end scheduler determines where to assign incoming queries considering how to maximize reuse of cached objects [1]. The index-based scheduling policies cause higher communication overhead on the front-end scheduler, and the cache index may not predict contents of the cache accurately if there are a large number of queries waiting to execute in the back-end application servers.

## III. DISTRIBUTED AND PARALLEL QUERY PROCESSING FRAMEWORK

Figure 1 shows how a scheduler in our distributed multiple query processing framework makes a scheduling decision for a given query. The goal of this framework is parallel processing of multiple queries, load balancing, and high ratio of cached data reuse via distributed semantic caching system. The front-end scheduler interacts with clients for receiving clients' queries and forwards the query to one of the back-end application servers. Each back-end application server independently runs on a cluster node and has an access to large-scale back-end storage devices or networked databases.

When a query is submitted to the front-end scheduler, it determines which back-end application server should process the query. Most schedulers in distributed query process-

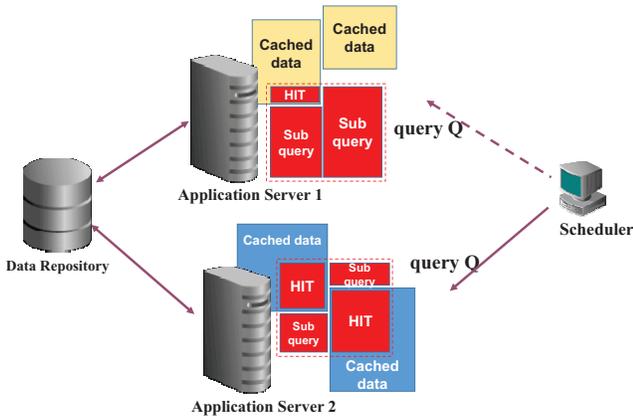


Fig. 1. Architecture of Distributed Query Processing Framework

ing middlewares employ a monitoring service that periodically collects performance metrics from back-end application servers. Based on the performance metrics, the schedulers make scheduling decisions to balance the system load across back-end applications.

Our distributed query processing framework is a component-based distributed job processing system that is currently under development. Our framework provides a programming model that scientific data analysis application developers can use to implement user-defined functions to process multi-dimensional scientific datasets on top of our framework. The programming model in our framework includes APIs that search and read raw datasets from back-end storage devices using multi-dimensional metadata of the dataset, interfaces that application developers can implement for application specific user-defined operations, and APIs that store and search intermediate query results.

When a query is forwarded to an application server, the server executes an application-specific user-defined operator for the query. The user-defined operator searches for cached objects in its semantic cache that can be reused to either completely or partially answer the query. Our framework first finds a cached data object that has the largest overlap with the query. If the reusable cached object doesn't cover the whole range of the query, then it partitions the query into several sub-queries for the partial region that was not in the cache as in Figure 1, and repeats the same process for the sub-queries recursively. If no cached object is found in the cache, it reads raw dataset from storage systems and process the sub-query from scratch.

Since reading raw dataset and generating intermediate query output is a very expensive operation for large scale scientific data analysis application, if possible, the scheduler needs to make scheduling decisions to maximize cached data reuse in distributed caching facility. For example, in Figure 1, a given query  $Q$  better be forwarded to application server 2 since its cached multi-dimensional data objects overlap a larger portion of the query than application server 1. However, if a scheduler makes scheduling decisions only based on cache hit ratio, it may cause system load imbalance and hurt overall system throughput. For example, if a very

few application servers have popular cached data objects, the application servers will be flooded with most of incoming queries, and it hurts system load balance and as a result system throughput will be poor.

#### IV. PREDICTING CACHED DATA OBJECTS IN DISTRIBUTED SEMANTIC CACHES

DEMA (Distributed Exponential Moving Average) that we propose is a multiple query scheduling policy that improves system throughput and reduces the query response time by taking into account of both load balance and cache hit ratio at the same time using a statistical prediction method. As distributed caching system dynamically changes its cached contents at very fast rate, DEMA scheduling policy makes real time scheduling decisions based on aggregated statistics of historical queries.

##### A. Background: Exponential Moving Average (EMA)

Our DEMA scheduling policy aggregates statistics of past queries using exponential moving average method. Exponential moving average (EMA) is a well-known statistical method to obtain long-term trends and smooth out short-term fluctuations, which is commonly used in stock price and trading volume predictions. EMA formula computes a weighted average of all historical data by assigning exponentially more weight to recent data. As our target applications process multi-dimensional range query, we compute multi-dimensional exponential moving average point as follows.

Let  $p_t$  be the cached object at time  $t > 0$  and  $EMA_t$  be the computed average at time  $t$  after adding  $p_t$  into the cache. Given the *smoothing factor*  $\alpha \in (0, 1)$  and the previous average  $EMA_{t-1}$ , the next average  $EMA_t$  can be defined incrementally by

$$EMA_t = \alpha \cdot p_t + (1 - \alpha) \cdot EMA_{t-1} \quad (1)$$

and it can be expanded as

$$EMA_t = \alpha p_t + \alpha(1 - \alpha)p_{t-1} + \alpha(1 - \alpha)^2 p_{t-2} + \dots \quad (2)$$

Since our framework is supposed to run long enough to amortize the initial EMA error, we simply chose  $k = 1$ , i.e.,  $EMA_1 = p_1$  where  $p_1$  is the first observed data. Equation 1 is used for the following EMAs. Note that EMA value can be a multi-dimensional element depending on the data domain of the application.

The smoothing factor  $\alpha$  determines the degree of weighing decay towards the past. For example,  $\alpha$  close to 1 drastically decreases the weight on the past data (short-tailed) and  $\alpha$  close to 0 gradually decreases (long-tailed).

##### B. Distributed Exponential Moving Average (DEMA)

The *Distributed Exponential Moving Average (DEMA)* scheduling policy employs as many EMA values as the back-end application servers to determine which server is to be assigned for each query and how the assignment adapts to the dynamic change of the distribution of input datasets. As our target applications process multi-dimensional datasets,

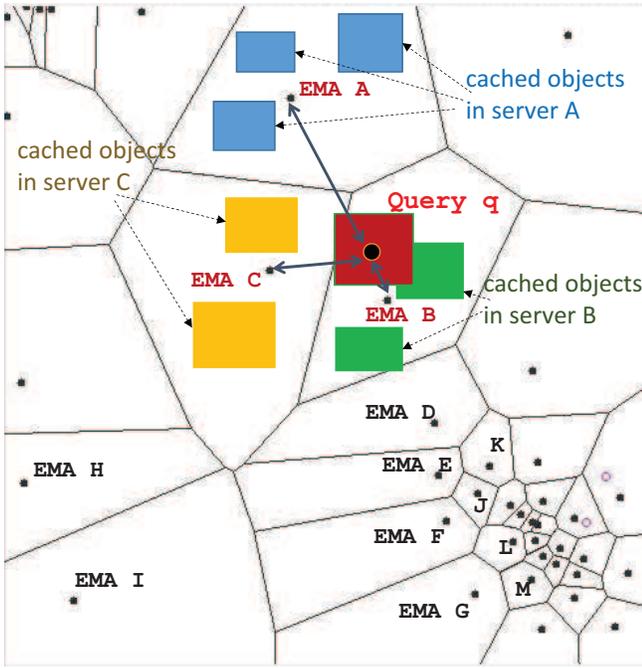


Fig. 2. DEMA scheduler calculates the Euclidean distance between EMA points and an input data needed by an incoming job, and assigns the job to the server  $B$  whose EMA point is closest.

cached objects in each back-end application server can be represented by multi-dimensional coordinates. As shown in Figure 2, DEMA scheduling policy computes a representative multi-dimensional point (*EMA point*) for the list of cached objects in each application server's cache. When a new query is submitted to a scheduler, the scheduler computes the distance between the query and EMA points. In the example shown in the Figure 2 the query is closer to server 1's cached objects. Hence the query should be forwarded to server 1, then the application server can reuse some of the cached objects in order to process the new query.

Our DEMA scheduling algorithm works in the following three steps.

a) 1. *Initialization.*: For simplicity, let a unit circle represent the *query space* where the value starts from zero and increases in a clockwise direction up to one which meets at zero. We assume that each input data is mapped to a point in multi-dimensional problem space using its center. Suppose there are  $n$  back-end application servers that process user jobs and one scheduling server that assigns each incoming job to an application server, as shown in Figure 1. In the beginning, the scheduling server creates  $n$  EMA variables  $EMA_1, EMA_2, \dots, EMA_n$ , with  $EMA_i$  for the  $i$ th application server, and then assigns the first  $n$  input data to  $EMA_1, EMA_2, \dots, EMA_n$ , respectively.

b) 2. *Assignment.*: Let  $q$  denote the center point of an input data requested by an incoming job, ranging from 0 to 1, and let  $q$  also represent the job itself. Upon receiving a job  $q$ , the scheduling server finds the EMA point  $EMA_{i^*}$  that is closest to  $q$  (in terms of their Euclidean distance), and then assigns  $q$  to server- $i^*$ . For example, in Figure 2 an input data  $q$  is closest to EMA point  $B$ , hence an incoming

job that needs the input data  $q$  is assigned to server  $B$ .

c) 3. *Update.*: Once  $q$  is assigned to server- $i^*$ , the scheduling server updates the EMA point for the server by

$$EMA_{i^*} = \alpha q + (1 - \alpha)EMA_{i^*}. \quad (3)$$

d) *Complexity.*: In step 2, the scheduling server needs to find the closest EMA point to a given input data. This algorithmic problem, known as Nearest Neighbor Search, can be solved in  $O(\log n)$  using space partitioning data structure such as *k-d tree* or *R-tree*.<sup>1</sup>

### C. Improving Cache Hit

Note that an application server keeps only certain number of recent input data in its cache; for example, it may employ Least Recently Used (LRU) policy. Let  $L$  denote the current number of cached input data in the cache. Ideally, it is desired that the EMA point reflects only those in the cache. However, EMA reflects all the past data including those that were expunged from the cache, but with less weight for older ones. Therefore, we may want to choose the decay factor  $\alpha$  such that the weight sum of the expunged input data is managed below a threshold. Formally, given a threshold  $\epsilon \in (0, 1)$ , we want

$$(1 - \alpha)^L < \epsilon. \quad (4)$$

Since  $L$  depends on the size of cached objects, we use the average number of input data in the cache instead, denoted by  $L^*$ . Therefore, we can choose the smallest  $\alpha$  that meets

$$\alpha > 1 - \epsilon^{L^*}. \quad (5)$$

In this way, the EMA value of an application server reflects what are current in the cache, and assigning incoming jobs that are close to the EMA, only jobs that need similar input data (i.e., close to each other in problem space) are assigned to the application server. This clustering effect promotes the cache hit ratio in the application server.

### D. Load Balancing of DEMA

Another important feature that distributed multiple query scheduling policies guarantee is *load balancing*; each job can be assigned to one of the application servers with equal probability, leveling the workloads of the application servers. To that end, DEMA scheduling algorithm gradually moves the EMA points so that they follow the distribution of requested input data. If a certain region of problem space becomes hot, the EMA points move to the hot spot by EMA equation. If input data requests are uniformly distributed, then DEMA algorithm tries to keep the sizes of Voronoi regions as similar as possible. DEMA algorithm achieves this goal by evenly distributing EMA points throughout the problem space.

Let us assume that the problem space is a 2-dimensional plane and the input data requests follow a uniform distribution. Figure 2 is a snapshot of the Voronoi diagram where a

<sup>1</sup>In our implementation, however, we performed a linear search that takes  $O(n)$  assuming that  $n$  is small enough, for example, up to 40.

dot represents an EMA value of an application server (e.g., EMA  $B$  is the EMA value of server  $B$ ), and a line segment represents the set of Euclidean middle points between two EMA points. In DEMA scheduling, an incoming query that falls in a Voronoi cell of an EMA point makes the EMA point slide toward the input data point, where the amount of the movement is affected by the smoothing factor  $\alpha$  (the more  $\alpha$ , the more shift). In Figure 2, the job that needs the input data  $q$  will be forwarded to server  $B$  since the EMA point of server  $B$  is closer to the query than any other EMA points. It is called the *Voronoi assignment model* where we assign every multidimensional point to the nearest cell point. The query assignment regions induced from the DEMA query assignment form a *Voronoi diagram* [14].

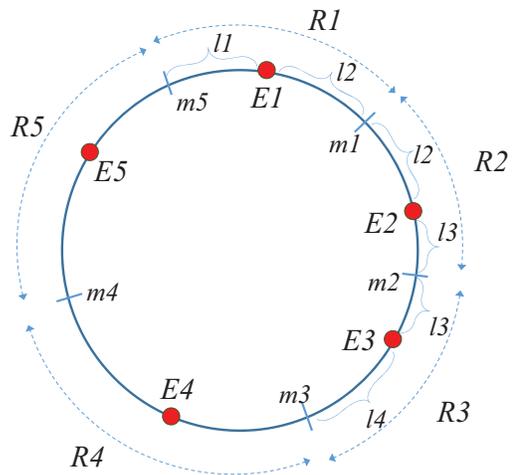
In this case, EMA  $B$  moves a little toward  $q$ , moving all the border lines of the cell a little toward the same direction. An important observation is that each EMA has tendency to move toward the center of its cell. For example, EMA  $E$  in the figure is located at the right corner of the cell. Thus more queries will arrive to the left of the EMA in the cell than to the right. Therefore, the EMA  $E$  is more likely to slide to the left rather than to the right. This movement will also move the border to the left, making the Voronoi region of EMA  $I$  smaller and EMA  $J$  and  $K$  larger. This property hints at the mechanism of load balancing; DEMA scheduling algorithm has a tendency to make a larger cell (EMA  $I$ ) smaller and a smaller cell (EMA  $K$ ) larger.

### E. Proof of DEMA Load Balancing

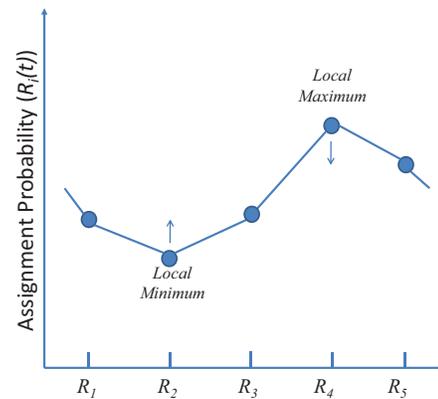
In this section, we provide a mathematical exposition of why DEMA balances the workload of application servers when the distribution of input data follows a uniform distribution. For brevity, we assume a 1-dimensional problem space in the proof. However, one can apply the same proof method described below to any higher dimension without difficulty.

Suppose we are given a 1-dimensional problem space as a unit circle (see Figure 3(a)) with five EMA values  $\{E_1, E_2, \dots, E_5\}$  representing application servers  $\{S_1, S_2, \dots, S_5\}$ , respectively. Let us denote the middle point of  $E_i$  and  $E_{i+1}$  by  $M_i$ . Then, the Voronoi cell region of  $S_i$  ( $i > 0$ ) is represented by a range  $[m_{i-1}, m_i]$  denoted by  $R_i$ . We use  $R_i$  for both the range and the length of the range interchangeably. In DEMA, all the jobs that need input data whose center falls in  $R_i$  are assigned to server  $S_i$ . Since the total size of the problem space is one and uniform query distribution is assumed,  $R_i$  also represents the probability that a query is assigned to server  $S_i$ . Note that  $\sum_i R_i = 1$ .

In fact, we can view the whole DEMA system as a Time Series; each  $E_i$  (thus  $R_i$ ) is a time series  $E_i(t)$  (thus  $R_i(t)$ ) that changes over a discrete time  $t = 1, 2, 3, \dots$  and the transition from  $t$  to  $t + 1$  happens when a job that needs data point  $q(t)$  is assigned to server  $i$ . In other words, given the current state  $\langle E_1(t), \dots, E_5(t) \rangle$ , the next job that needs a data point  $q(t)$  triggers a state transition to  $\langle E_1(t + 1), \dots, E_5(t + 1) \rangle$ . Since the transition only depends on the previous state, the whole DEMA system



(a) Problem Space with five EMA values at time  $t$



(b) Load Balancing by Convex Optimization

Fig. 3. DEMA Load Balancing in One-dimensional Space

is also a Markov Chain with discrete time on a general state space [15]. As described below, however we take a simpler approach than a Markov chain analysis to prove load balancing.

Figure 3(b) depicts the sizes of the ranges for all the servers corresponding to Figure 3(a) at time  $t$ . In the figure, we note that  $R_2(t)$  is a local minimum and  $R_4(t)$  is a local maximum, that is,  $R_2(t) < R_1(t) \wedge R_2(t) < R_3(t)$ , and  $R_4(t) > R_3(t)$ ,  $R_4(t) > R_5(t)$ . To balance the workloads of the application servers, DEMA algorithm tries to make the assignment probability of all the servers as similar as possible. It does so by, in the next step, decreasing all the local maxima and increasing all the local minima. Formally, we want

$$E[R_2(t + 1)] > R_2(t) \quad (6)$$

and

$$E[R_4(t + 1)] < R_4(t) \quad (7)$$

where  $E[\cdot]$  is the expectation function. With this property, the system is intuitively moving toward a better balanced state.

*Theorem 1:* For a local maximum  $R_i(t)$  at time  $t$ , we have

$$E[R_i(t+1)] < R_i(t) \quad (8)$$

and for a local minimum  $R_i(t)$  at time  $t$ , we have

$$E[R_i(t+1)] > R_i(t). \quad (9)$$

*Proof:* Due to similarity, we only prove Equation 9. In this proof, we directly use the example in Figure 3(a). However, one can easily change the proof to a general case. In the following, when there is no confusion, we omit the time index  $t$  for simplicity (e.g.,  $E_2(t) = E_2$ ). In Figure 3(a)  $l_i$  is the distance from  $m_{i-1}$  to  $E_i$  and the distance from  $E_{i-1}$  to  $m_{i-1}$ . Without loss of generality, we assume that  $m_5 = m_0 = 0$ . Note that we have a local minimum  $R_2(t)$ ;  $R_2(t) < R_1(t)$  and  $R_2(t) < R_3(t)$ , or equivalently

$$\ell_3 < \ell_1 \text{ and} \quad (10)$$

$$\ell_2 < \ell_4. \quad (11)$$

Suppose a new job that needs input data  $x \in [0, 1)$  arrives. If  $x$  falls outside  $\overline{m_5 m_3}$ , none of  $E_1, E_2$ , and  $E_3$  moves, thus we get  $R_2(t+1) = R_2(t)$ . If  $x \in [0, m_1]$  then  $E_1(t+1) = \alpha x + (1-\alpha)E_1(t)$ , thus we get

$$\begin{aligned} m_1(t+1) &= \frac{1}{2}(E_1(t+1) + E_2(t+1)) \quad (12) \\ &= \frac{1}{2}(\alpha x + (1-\alpha)E_1(t) + E_2(t)) \quad (13) \end{aligned}$$

If  $x \in [m_1, m_2]$ ,  $E_2$  moves but the length of  $R_2$  remains the same, i.e.,  $R_2(t+1) = R_2(t)$ . When  $x \in [m_2, m_3]$ , then  $E_3(t+1) = \alpha x + (1-\alpha)E_3(t)$ , thus we get

$$\begin{aligned} m_2(t+1) &= \frac{1}{2}(E_2(t+1) + E_3(t+1)) \quad (14) \\ &= \frac{1}{2}(E_2(t) + \alpha x + (1-\alpha)E_3(t)). \quad (15) \end{aligned}$$

Now let us compute  $E[R_2(t+1)]$ . Since input data  $x$  follows a uniform distribution in  $[0, 1)$ , the probability density function of  $x$  is  $f(x) = 1$ , and the cumulative density function is  $F(x) = x$ . By definition,

$$\begin{aligned} E[R_2(t+1)] &= \int_0^{m_1} m_2 - \frac{1}{2}(\alpha x + (1-\alpha)E_1 + E_2) dx \\ &+ \int_{m_1}^{m_2} R_2 dx \\ &+ \int_{m_2}^{m_3} \frac{1}{2}(E_2 + \alpha x + (1-\alpha)E_3) - m_1 dx \\ &+ \int_{m_3}^1 R_2 dx \quad (16) \end{aligned}$$

Since  $E_1 = \ell_1, E_2 = \ell_1 + 2\ell_2, E_3 = \ell_1 + 2\ell_2 + 2\ell_3$  and  $m_i = E_i + \ell_{i+1}$  for  $i = 1, 2, 3$ , the average change of  $R_2$  becomes

$$\begin{aligned} E[R_2(t+1)] - R_2(t) &= E[R_2(t+1)] - (\ell_2 + \ell_3) \\ &= \frac{\alpha}{4}(\ell_1^2 - \ell_3^2 + \ell_4^2 - \ell_2^2). \quad (17) \end{aligned}$$

Because of Equation 10 and 11, we have

$$E[R_2(t+1)] - R_2(t) > 0. \quad (18)$$

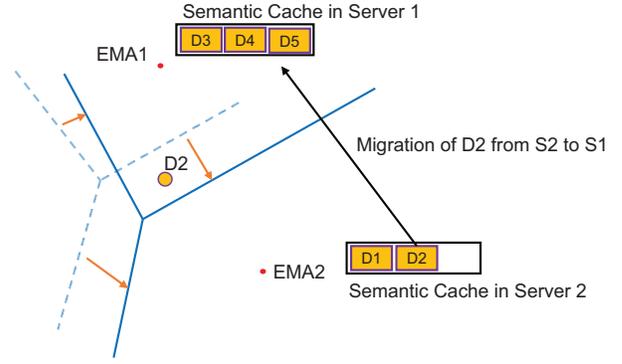


Fig. 4. An Example of Data Migration: Location of a cached data object  $D2$  is not covered by server  $S2$ 's region since the boundaries change. If  $D2$  is a frequently accessed data object, it better be moved to server  $S1$ .

Therefore, on average, the assignment probability of a local minimum server is more likely to increase at the next query. ■

It is possible that decreasing a local maximum can make a neighboring interval a new local maximum (and the same applies to local minimum). Nevertheless, subsequent jobs are likely to decrease those new local maxima (or increase new local minima) as shown above. However, this only happens when the assignment probabilities of the local maxima and neighboring EMAs are close to each other, and in that case it is likely that the overall balance was already achieved.

## V. CACHED DATA MIGRATION

As various incoming queries access different parts of input datasets, EMA points dynamically move to new locations over time. As a result of the movement, some cached data objects may cross the boundaries of Voronoi regions. If the changed Voronoi region does not contain a cached data object any more, the cached data object is not likely to be used by subsequent queries because DEMA scheduler will not forward a query to the server.

Figure 4 shows an example of Voronoi region changes. If the boundary of Voronoi region moves to a lower right direction, some of the cached data objects in server  $S2$  will not be covered by the updated  $S2$ 's region. We will refer to the cached data objects out of current Voronoi region as "misplaced cached data". In the example, a data object  $D2$  is a misplaced cached data. Based on the updated boundaries, the scheduler will not forward any incoming query that needs data  $D2$  to server  $S2$ , but instead it will forward it to server  $S1$ . Although a previously computed query result  $D2$  is in a neighbor server  $S2$ , server  $S1$  will read raw datasets from storage systems and process the query from scratch. Since  $S2$  is not likely to reuse data object  $D2$ ,  $D2$  will be eventually evicted from  $S2$ 's cache.

In order to manage the distributed semantic buffer cache seamlessly, we propose data migration of cached data objects to improve overall cache hit ratio. With the data migration policy enabled, the misplaced cached data objects are mi-

grated to a remote server whose Voronoi region encloses the cached data objects.

When a front-end server receives a query, it searches for a back-end application server whose EMA point is the closest to the query point. Periodically or when EMA points have moved by a large margin, the front-end server constructs a piggyback message with the server's neighbor EMA points and forwards the query and the neighbor EMA points to the selected back-end application server. With the updated EMA points, back-end application servers can migrate misplaced cached data objects in either pull mode or push mode.

1) *Pull Mode Migration*: When a back-end application server searches for cached data objects in its own cache but does not find it, pull mode migration policy searches for the cached data objects in its neighbor servers' cache. Since it is very expensive operation to look up neighbor servers' cache for every cache miss, each back-end application server in our framework periodically collects EMA points of its neighbor servers. If a cache miss occurs in a server  $s$ , the server checks if the query's point is closer to a neighbor server's historical EMA point rather than its own previous EMA point. If another server  $n$ 's previous EMA point is closer to the query than its previous EMA point, it means the query would have been assigned to the neighbor server  $n$  unless the EMA points moved, and it is highly likely that the neighbor server  $n$  has a misplaced cached object for the query. Thus the application server  $s$  should check the neighbor server  $n$  and pulls the cached data object from  $n$ .

The performance of pull mode migration depends on how accurately each application server identifies a neighbor server that has misplaced cached data objects. Searching all neighbor servers for every cache miss can increase cache hit ratio but it will place too much overhead on the cache look up operation, hence we implemented a push mode migration as described below.

2) *Push Mode Migration*: When a back-end server receives a query, first it searches for its own cached objects in its semantic cache that can be reused. During the look up operation, push mode migration policy investigates whether each cached data object is misplaced or not. If a cached data object is closer to a neighbor EMA point than its own EMA point, the cached data object is considered to be misplaced and we migrate it to a neighbor server. When the neighbor server receives a migrated data object, it determines if the migrated data object should be stored in its cache using its cache replacement policy. Note that our data migration policy does not examine the entire cached data objects since it will cause significant overhead. Instead, we detect misplaced data objects only during cached data look-up operation and transfer them to neighbor servers. I.e., if we use a hash table for cached data look-up, only the cached data objects in an accessed hash bucket will be considered for data migration candidates. This policy migrates some misplaced cached data objects to neighbor servers in a lazy manner, which helps reducing the overhead of cache look-up operation.

As we will show in section VI, our experimental study shows data migration does not always guarantee higher cache hit ratio if the distribution of arriving queries is rather static.

With static query distribution, the boundaries of Voronoi regions may fluctuate but not by a large margin. In such cases, cached data objects near boundaries tend to be migrated between two servers repeatedly, and it is possible for them to evict some useful cached data objects continuously in neighbor servers. There can be various solutions to mitigate the negative effect of the continuous eviction problem. One is to keep cached data objects unless its center point is not significantly close to a neighbor EMA point. Another option is to replicate the cached data objects on the boundaries by not deleting migrated cached data objects and we mark recently migrated data objects so that the same data objects are not migrated again. In our implementation we chose the second option, which prevents unnecessary repeated LRU cache updates.

## VI. EXPERIMENTS

### A. Query Workload Generator

In order to show that the DEMA scheduling policies perform well with various input data distributions in multi-dimensional space, we synthetically generated 40,000 2-dimensional queries using uniform, normal, and Zipf's distributions respectively. For all the workloads, the query inter-arrival time distribution was modeled by Poisson process. Since frequently accessed hot data objects in these synthetic workloads do not change over time, we also generated dynamic query workloads that change the requested input data distribution along time dimension. The dynamic query distribution mainly consists of normal distributions, but we made the mean and the standard deviation change every interval of 10,000 queries, i.e., the dynamic query distribution is composed of four successive normal distributions with different mean point and different standard deviation containing 10,000 queries each.

In addition to the synthetic probability distributions, we generated more realistic query workloads using a probability model - Customer Behavior Model Graph (CBMG). In this model, there are a large number of hot multi-dimensional regions of interest, and the first query is selected from one of the regions. Subsequent queries after the first one in the batch may either remain around that point (moving around its neighborhood, temporal movement, or resolution increase or decrease) or move on to a new hot region. This workload model is known to resemble real users' query access patterns and is commonly used to evaluate e-business applications and also used for web site capacity planning [16]. Using the CBMG model, we generated 40,000 2-dimensional queries (latitude, longitude) with various transition probabilities, but we only show the results using one of the transition probabilities because results for other transition probabilities were similar to those shown.

### B. Experiments

1) *Scalability*: Figure 5 depicts the average query response time of 40,000 queries and cache hit ratio for DEMA scheduling (*DEMA*), DEMA scheduling with data migration (*MDEMA*), and load-based scheduling (*LOAD*), as

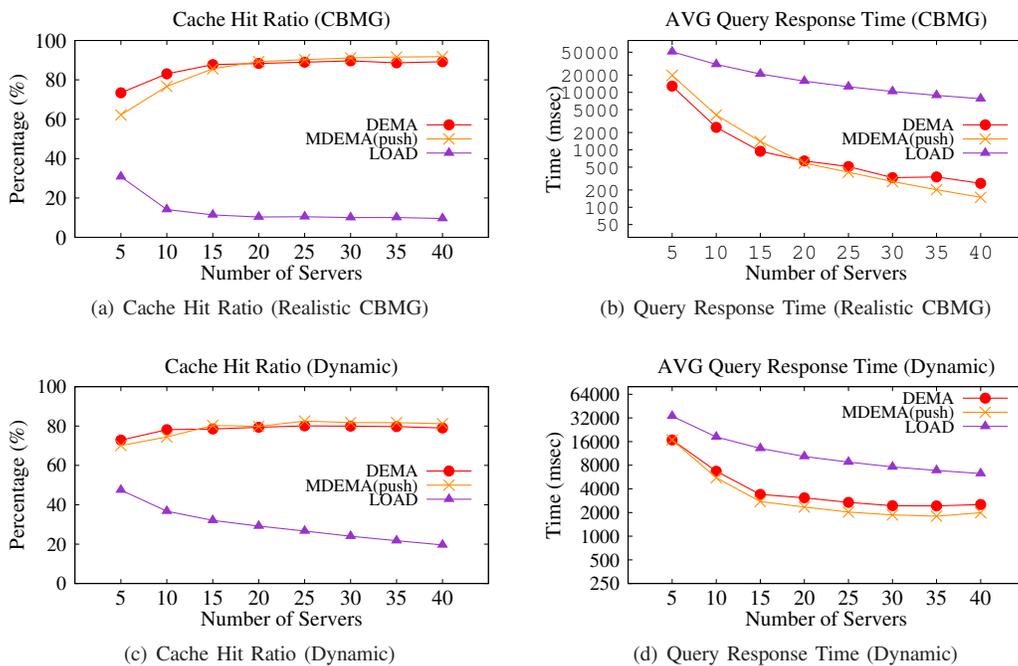


Fig. 5. Performance Comparison Varying Number of Servers

the number of back-end application servers increases. The query response time is defined as the amount of time from the moment a query is submitted to the scheduler until it completes in a back-end application server, i.e., it includes the waiting time in the application server’s queue as well as the actual processing time. The query inter-arrival time was modeled by Poisson Process with average of 1 ms, and each back-end application server has a cache size that can hold up to 200 query results. For the rest of the experiments, we fixed the EMA smoothing factor  $\alpha$  to 0.03. The smoothing factor determines how much weight is given to each past query result. When  $\alpha$  is 0.03, the recent 200 query results account for 99.7% of the weight in the EMA equation.

As the number of servers increases, DEMA scheduling policy improves the cache hit ratio since the total available cache size in distributed caching system increases. However, Load-based scheduling policy does not get any benefit of the increased cache size but its cache hit ratio decreases because frequently requested cached data objects are scattered across more back-end application servers and load-based scheduling policy does not know which remote cache has what cached data objects. Hence the performance gap between load-based scheduling policy and DEMA scheduling policy becomes wider as the number of servers increases. Due to its low cache hit ratio, the average query response time of load-based scheduling policy is order of magnitude higher than that of DEMA scheduling policy.

With a small number of servers, the total size of distributed caching system is not enough to hold all frequently accessed data objects, hence the migration of misplaced cached objects evicts other frequently accessed data objects in neighbor servers. It results in lowering overall cache hit ratio in our experiments, thus DEMA scheduling policy

that does not migrate misplaced cached objects outperforms MDEMA. However with a larger number of servers, the total size of distributed caching system becomes large enough to hold most of the frequently accessed data objects. If each application server’s cache is large enough, some of the cached data objects are less likely to be accessed than others, thus our cache migration policy evicts such less popular data objects and replaces them with more popular data objects from neighbor servers, which are highly likely to be accessed by subsequent queries. Thus with more than 25 back-end application servers, DEMA with data migration policy outperforms DEMA scheduling policy.

2) *Cache Miss Overhead:* Figure 6 depicts the average query response time for uniform, normal, Zipf’s, and dynamic input data distributions with varying the cache miss penalty. If a query needs a larger raw data object and more computation, cache miss causes a higher performance penalty. For the experiments, we fixed the cache size so that it can hold at most 200 query results, and employed 40 back-end application servers. In the Figure 6, we do not present the cache hit ratio and the standard deviation of server loads since the effect of cache miss penalty is not relevant to the cache hit ratio and system load balance. As the cache miss overhead increases, the average query response time linearly increases not only because of the cache miss penalty but also because the waiting time in the queue increases. The waiting time increases because previously scheduled queries suffer from higher cache miss penalty and more queries will be accumulated in the queue. In the experiments, average query response time with load-based scheduling policy is up to 72 times higher than DEMA scheduling policy. When data migration is employed with DEMA scheduling policy the hit ratio improves from 57~83% to 55~92%, thus it

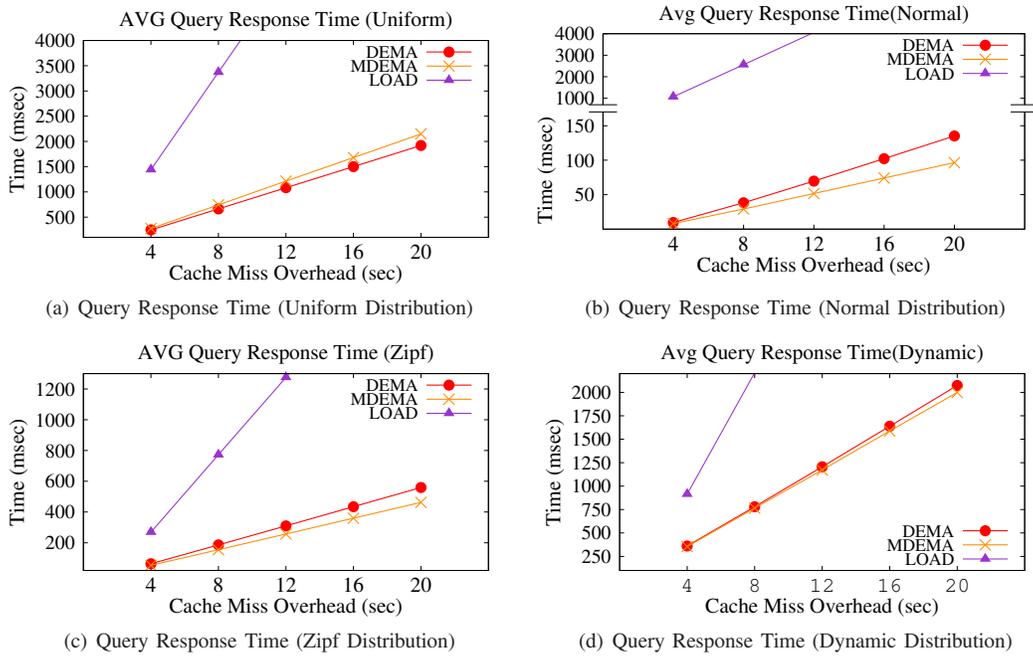


Fig. 6. The Effect of Cache Miss Overhead

exhibits up to 40% lower average query response time. Load-based scheduling policy exhibits significantly higher average query response time because of low cache hit ratio. Hence note that y-axis in Figure 6(b) is broken into two ranges to show the huge performance gap between load-based scheduling policy and DEMA scheduling policy. The graphs in Figure 6 show the average query response time of DEMA and MDEMA is mostly smaller than the cache miss penalty, i.e., most queries reuse the cached data objects in distributed caching system and they are executed immediately without being accumulated in the waiting queue. But load-based scheduling policy suffers from low cache hit ratio and the cache miss penalty is higher than query inter-arrival time, thus its average query response time is much higher than the query inter-arrival time because the queries wait in the queue for significant amount of time.

3) *Cache Size Effect*: Now we consider the impact of cache size to the system. We ran experiments with 40 back-end application servers, and fixed the EMA smoothing factor to 0.03 for DEMA policy. The cache miss penalty was 20 ms. In the graphs shown in Figure 7, Load-based scheduling policy gains the benefit of increased cache size. Note that load-based scheduling policy didn't get benefit of increased number of distributed caches. With non-intelligent scheduling policies such as round-robin or load-based scheduling policy, each application server should have very large cache in order to exhibit high cache hit ratio.

Unlike load-based scheduling policy, DEMA scheduling policy consistently shows very high cache hit ratio (> 80%) even when cache size is small. This is because DEMA scheduling policy considers all cached objects in distributed caching system. Thus even if a single server's cache size is small, the total size of distributed caching system is not small. And DEMA scheduling policy assigns each query to

a back-end application server which is highly likely to have the requested data object in its cache even when the cache size is small.

It should be noted that the cache hit ratio is totally dependent on query workloads. If certain query workloads do not have query sub-expression commonality, even DEMA scheduling policy may have low cache hit ratio. However, load-based scheduling policies will have even lower cache hit ratio than DEMA. If cache hit ratio is low, a large number of queries have to be computed from the scratch. In such situations, load balancing plays an important role in decreasing the wait time in the queue. DEMA scheduling policy shows higher standard deviation of the number of processed queries in each server than load-based scheduling policy, but since cache hit does not put significantly high overhead to back-end application server, we counted only the number of missed queries in each server and DEMA scheduling policy exhibited very good load balancing behavior.

In Figure 7, the benefit of migrating misplaced cached objects is not clearly observable when cache sizes are small. Again this is because the migration may evict other frequently accessed data objects in neighbor servers. However, as the cache size increases, data migration improves the cache hit ratio and it results in faster query response time.

## VII. CONCLUSION

This paper presents a multiple query scheduling policy for distributed query processing framework that takes into consideration the dynamic contents of distributed caching infrastructure. In distributed query processing systems where the caching infrastructure is distributed and scales with the number of servers, both leveraging cached results and achieving load balance become equally important to improve the overall system throughput. In order to achieve load

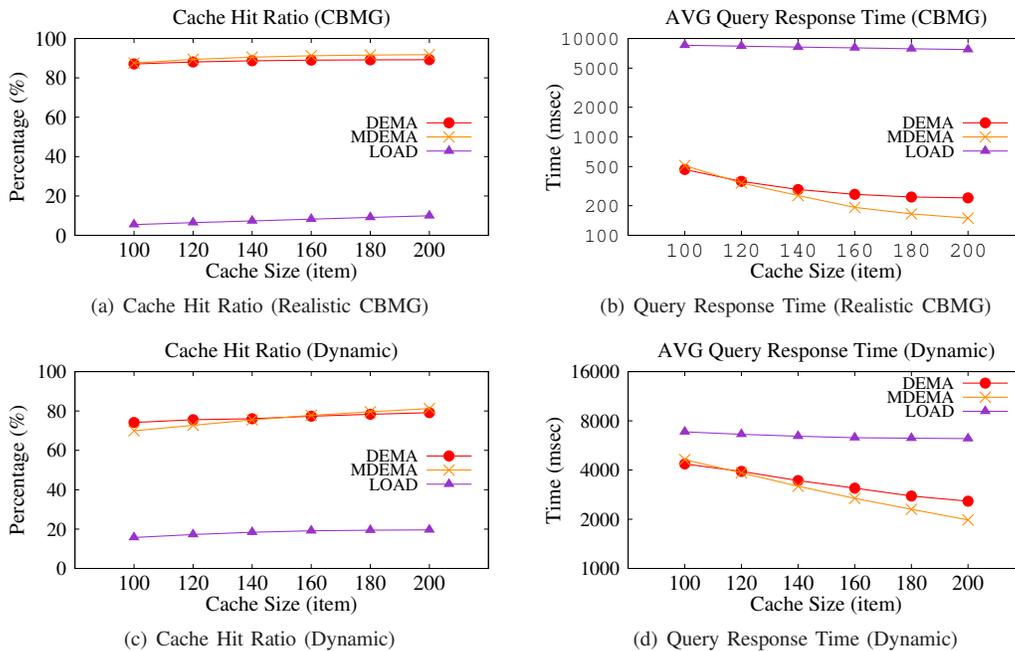


Fig. 7. Effect of Cache Size with Realistic CBMG Workload and Dynamic Workload

balancing as well as to exploit cached query results, it is required to employ more intelligent query scheduling policies than the traditional round-robin and load-based scheduling policies.

Leveraging distributed cache data is very important for improving system performance when applications can take the advantage of sub-expression commonality of cached results. Our framework implemented a cached data object migration module that helps distributed caches cooperate to further improve cache hit ratio. In our experiments, we showed our cache-aware query scheduling policy with data migration policy outperforms legacy load-based scheduling policy.

#### REFERENCES

- [1] B. Nam, M. Shin, H. Andrade, and A. Sussman, "Multiple query scheduling for distributed semantic caches," *Journal of Parallel and Distributed Computing*, vol. 70, no. 5, pp. 598–611, 2010.
- [2] B. Godfrey, K. Lakshminarayanan, S. Surana, R. Karp, and I. Stoica, "Load balancing in dynamic structured p2p systems," in *Proceedings of INFOCOM 2004*, 2004.
- [3] U. V. Catalyurek, E. G. Boman, K. D. Devine, D. Bozdag, R. T. Heaphy, and L. A. Riesen, "A repartitioning hypergraph model for dynamic load balancing," *Journal of Parallel and Distributed Computing*, vol. 69, no. 8, pp. 711–724, 2009.
- [4] N. Vydyanathan, S. Krishnamoorthy, G. Sabin, U. Catalyurek, T. Kurc, P. Sadayappan, and J. Saltz, "An integrated approach to locality-conscious processor allocation and scheduling of mixed-parallel applications," *IEEE Transactions on Parallel and Distributed Systems*, vol. 15, pp. 3319–3332, 2009.
- [5] Q. Zhang, A. Riska, W. Sun, E. Smiri, and G. Ciardo, "Workload-aware load balancing for clustered web servers," *IEEE Transactions on Parallel and Distributed Systems*, vol. 16, no. 3, pp. 219–233, 2005.
- [6] J. L. Wolf and P. S. Yu, "Load balancing for clustered web farms," *ACM SIGMETRICS Performance Evaluation Review*, vol. 28, no. 4, pp. 11–13, 2001.
- [7] M. Katevenis, S. Sidiropoulos, and C. Courcoubetis, "Weighted round-robin cell multiplexing in a general-purpose atm switch chip," *IEEE Journal on Selected Areas in Communications*, vol. 9, no. 8, pp. 1265–1279, 1991.
- [8] M. Aron, D. Sanders, P. Druschel, and W. Zwaenepoel, "Scalable content-aware request distribution in cluster-based network servers," in *Proceedings of Usenix Annual Technical Conference*, 2000.
- [9] V. Pai, M. Aron, G. Banga, M. Svendsen, P. Druschel, W. Zwaenepoel, and E. Nahum, "Locality-aware request distribution in cluster-based network servers," in *Proceedings of ACM ASPLOS*, 1998.
- [10] J.-S. Kim, H. Andrade, and A. Sussman, "Principles for designing data/compute-intensive distributed applications and middleware systems for heterogeneous environments," *Journal of Parallel and Distributed Computing*, vol. 67, no. 7, pp. 755–771, 2007.
- [11] K. Zhang, H. Andrade, L. Raschid, and A. Sussman, "Query planning for the Grid: Adapting to dynamic resource availability," in *Proceedings of the 5th IEEE/ACM International Symposium on Cluster Computing and the Grid (CCGrid)*, Cardiff, UK, May 2005.
- [12] M. Rodríguez-Martínez and N. Roussopoulos, "MOCHA: A self-extensible database middleware system for distributed data sources," in *Proceedings of 2000 ACM SIGMOD*.
- [13] J. Smith, S. Sampaio, P. Watson, and N. Paton, "The polar parallel object database server," *Distributed and Parallel Databases*, vol. 16, no. 3, pp. 275–319, 2004.
- [14] M. de Berg, O. Cheong, M. van Kreveld, and M. Overmars, *Computational Geometry, Algorithms and Applications*. Springer, 1998.
- [15] C. P. Robert and G. Casella, *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2005.
- [16] D. A. Menasce and V. A. F. Almeida, *Scaling for E-Business: Technologies, Models, Performance, and Capacity Planning*. Prentice Hall PTR, 2000.