

**Protecting the Quality of Service
of Existing Information Systems**

Kevin Beyer
Miron Livny
Raghu Ramakrishnan

Technical Report #1358

December 1997

Protecting the Quality of Service of Existing Information Systems

Kevin Beyer*

Miron Livny

Raghu Ramakrishnan

Computer Sciences Department

University of Wisconsin, Madison WI 53706

beyer, miron, raghu@cs.wisc.edu

Abstract

Most organizations that offer external access to their data would benefit from some mechanism that ensures a desired level of service for local users. In this paper, we propose such a mechanism, called the provider agent (PA) architecture, that protects local users by ensuring a (DBA specified) quality of service for local requests in the face of computational demands made by external requests. The PA is a general purpose solution that enhances most information systems currently available. The novelty of our approach is the combination of request profiling with load control mechanisms to improve both protection and performance, while not requiring any modifications to the underlying information system. We demonstrate the effectiveness of the proposed techniques with a prototype PA for a commercial DBMS.

1 Introduction

People want access to information within and across organizational boundaries. Many new tools are now available to allow people to cross old barriers, for example, Web browsers, CGI servers, Java¹ and Java Database Connectivity (JDBC), Microsoft's ODBC, IBM DataJoiner², and data warehouses. People are excited about the potential of these products, but what happens when a site decides to take advantage of this new technology and allow external users to access their existing database? The local users of the database have come to expect a certain level of performance. Suppose the site decides that it is willing to accept a controlled degradation in the performance of local requests to allow this new service, but the day-to-day business performed by the local users is the ultimate priority. If the external workload ever degrades the performance to an unacceptable level, the local users' goodwill might run out, meaning the elimination of external access.

*Contact Author: 608-262-6629

¹Java and Java-based marks are trademarks or registered trademarks of Sun Microsystems, Inc. in the United States and other countries.

²DataJoiner is a trademark of IBM Corporation.

In this environment, the local users are recognized as the owners of the system, and the external requests should be processed only when the owners allow it. If the site does not take precautions, that fateful day will arrive when the external workload brings the system to its knees, and correcting the problem could be far from trivial. The damage caused by a heavy external workload is usually a large increase in local response time, but the damage could also be unforeseen failures of local requests due to limited resources, for example, a lack of memory or connections.

One way to keep external activity from interfering with local users is to replicate the existing system. Local users continue to access the original system, while the replicated system services the external requests. The external requests cannot interfere with local users since they are not using the same system. However, this solution is not always acceptable for the following reasons. First, the cost of creating and administering a second system could be prohibitive, especially for large databases. Second, if the external users need access to up-to-date information, a replicated system is unacceptable. Third, a demand is still placed on the local users system when the replicated system is periodically updated (e.g., data warehouse updates), and that demand, although more predictable, still needs to be scheduled. Fourth, if the external requests are prioritized, replicating the system for each priority level is infeasible. Finally, replicating the system effectively partitions the computing resources, so one system could be idle while the other is overloaded. For these reasons, we do not consider the replicated solution further.

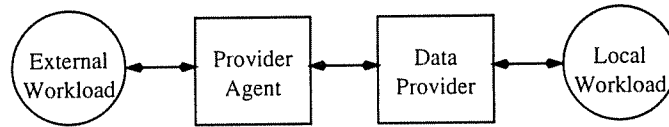


Figure 1: Provider Agent Architecture

This paper describes the architecture of a system, called the Provider Agent (PA), that protects the local quality of service of a data provider (i.e., any data producing system) in the face of external processing demands. The PA mediates between the provider and the external world (see Figure 1), much like most multi-database (MDBS) solutions [5, 26]. The architecture is designed to work with any type of data provider available today, for example, a DBMS, a Web server, an FTP server, or even data generating programs like a simulation system. The PA enhances these existing systems, complements current connectivity tools, and can be deployed without waiting for new versions of these systems. The PA does not mandate any changes to the underlying data provider or the existing (local) applications.

In designing the PA, we used dynamic load control techniques to protect the local users and request profiling to improve external response. Load control transforms the worst case external workload from disastrous to predictable and acceptable. The synergy between load control and profiling enables the administrator to specify load controls based upon query resource requirements rather than treating all requests the same, independent of the potential impact of the request. The combination of load control and profiling to protect the quality of service, while not requiring any modifications to the underlying system, makes the PA an attractive and novel solution.

To evaluate the effectiveness of these techniques, we implemented a prototype PA for a relational DBMS and cre-

ated a testbed for experimentation. We found that our load control mechanisms, which are implemented outside the DBMS, effectively maintain the local quality of service at a predefined level, and our profiling technique dramatically improves the throughput of mixed external workloads.

In the remainder of this paper we further motivate the PA by describing some potential applications in Section 2. Then, in Section 3 we discuss our measure of success. In Section 4 we detail the PA load control and profiling mechanisms. Section 5 describes an implementation of the PA for a commercial DBMS and Section 6 presents the performance results. The remaining sections describe related work (Section 7), future work (Section 8), and conclusions (Section 9).

2 Applications

The PA is applicable to a broad range of scenarios. For example, when an enterprise connects its disparate departmental databases into a multi-database, the PA would ensure that the large ad-hoc queries posed to the multi-database by someone in the marketing department do not overwhelm any of the underlying production databases. When a company offers order tracking information to Internet users, the PA prevents the status requests from disrupting order processing. Many other scenarios exist, and in this section, we describe a few more sample applications to further motivate the utility of the PA.

NASA's EOSDIS project is an excellent example of a large scale data provider that will be available to an enormous group of users. The Earth Observing System (EOS) satellites are expected to return 300 terabytes of data per year which will then be augmented with additional information [24]. This data will be stored in EOSDIS (EOS Data and Information System) and made available to anyone from NASA scientists to elementary school children. We believe that the PA architecture could be used to schedule requests from different groups of users. For example, NASA scientists could access EOSDIS directly, while anonymous Internet users and university scientists could access EOSDIS through a PA. The PA would ensure that the NASA scientists' quality of service needs are met by controlling the demands placed upon EOSDIS by the external requests. The PA would also distinguish between the university researchers and the anonymous users by allowing researchers' requests to be serviced more quickly than the anonymous requests.

On a smaller scale, consider the processing of data-mining queries or the extraction of production data into a data warehouse from a standard OLTP system. The response time of the data-mining queries is probably less critical than that of the standard business transactions like order processing. If the database system has a predictable local workload, the DBA could avoid some problems by prohibiting any of the less important queries from running during peak hours. This solution is not generally acceptable for a number of reasons. First, this all-or-nothing approach is extreme because it might be possible to access the system in a controlled manner during peak hours without dramatically affecting the OLTP queries, especially since some of the data-mining queries may be easy to answer. Second, the DBA must translate peak periods into specific times of the day, implying that the system will

not dynamically respond to changing workloads. Third, some coordination between the database site and the user issuing the request must exist to let the user know when they are allowed to submit their requests. In short, forcing the DBA to schedule transactions is not a flexible solution and it requires too much manual intervention.

As a final example, consider a Web server that is about to be searched by an automatic indexing robot. If the robot arrives during a peak period of the day, it would be useful to delay the robot's requests until an off-peak time or at least control the rate at which the robot accesses the data. The service provided by the indexing sites is invaluable to Web users that are trying to locate information, so simply rejecting the request is probably not acceptable. The current solution is to rely upon the robot writers to space requests to an individual server so that the robot does not flood the server [16, 11]. The first problem with this approach is that the Web site administrator must rely on the robot writer to be conscientious and careful when writing the robot. The second problem is that at different times during the day, the server might be able to process the robots requests at a much higher rate than others. Why should the robot make arbitrary decisions about when and how frequently the Web server is accessed when the server is much more capable to answer these questions? If the indexing robot were to access the Web server through a PA, then the PA would ensure that the robot's requests do not drastically interfere with the Web servers response to normal requests.

3 Quality of Service

The primary concern in this paper is the quality of service (QoS) offered to the local users. Although many QoS measures exist, we use changes in the response times experienced by the local users to illustrate the protection offered by a particular PA policy. Instead of summarizing the response times into a single number like the average response time, we use the distribution function³ of the local response times as our QoS measure. The distribution captures the performance that the local user can expect at all percentiles, rather than at just the average.

Only when the local users are satisfied can the PA consider improvements in external processing. Two scheduling policies cannot be compared on the basis of external performance without considering the local QoS. In comparing external performance, we use the throughput rather than response time because we are more interested in the amount of work completed than how quickly the work is performed.

4 Mechanisms

The PA achieves its objectives through a variety of load control mechanisms. Each of the mechanisms presented in this section describes a way that the PA can control when and how a particular external job is run. The merits of load control are well known, dating back to operating systems work of the '60s and '70s [14, 9]. The novelty of our approach is in extending these ideas to the protection of the local QoS in the presence of an unknown external

³The function, $F_X(x) \equiv P(X \leq x)$, is the probability that the random variable X is less than x .

workload without changing the existing system, in our definition of fractional MPL, and in combining load control with query profiling. This section describes our load control mechanisms in detail, after describing some load control basics.

When considering a given (fixed) system, the response time of a request is determined by its execution time, the time the job spends waiting for resources (e.g., waiting to perform a disk I/O or to acquire a lock), and the time that is “wasted” due to context switches and resource sharing (e.g., decreased buffer hits, increased number of disk seeks). The execution time of a job is determined by the system, so additional jobs and the PA cannot affect it. The goal of the PA is to find ways of limiting the increase in response time caused by the latter two sources of delays.

The time spent waiting for resources is proportional to the number of jobs concurrently accessing the system — as more jobs are added to the system, the waiting time increases. Wasted time is more difficult to characterize. The maximum amount of wasted time could be achieved when only two jobs run concurrently. Even when only one job is run at a time, the system might still waste time due to inter-job sharing. For example, job *A* scans a relation that fits exactly within the buffer pool, and job *B* scans another relation. Alternately running jobs *A* and *B* implies that job *A* must read the relation from disk each time it executes, but if job *B* were not run, then job *A* would only read the relation once. By decreasing the frequency of the *B* jobs, the amount of time spent re-reading *A*’s relation can also be decreased.

By controlling the maximum number of concurrent jobs and the frequency of the jobs, the sources of delays should in turn be controlled. This is the basic PA protection mechanism — to control the number of concurrent external jobs. But by running external jobs the number of local requests in the system will increase, so the PA must account not only for the delays caused directly by the external jobs, but also for the increase in local response times caused indirectly by the increased number of local jobs. Sections 4.1 through 4.4 describe the PA load control mechanisms.

4.1 Integer MPL

Adjusting the multiprogramming level (MPL), by definition, controls the maximum number of active requests in a system. By limiting the number of external jobs to *E*, the DBA controls the worst case performance of a the local workload. Theoretically, by constraining the external workload to any *E*, a stable local workload will never become unstable by running external jobs. This can be seen by looking at the effective service rate. The effective local service rate of any one service center is⁴:

$$\frac{1}{\mu_{eff}} = \frac{1}{\mu_{act}} * \frac{L}{L + E}$$

where *L* and *E* are the number of local and external requests respectively, and $1/\mu_{act}$ is the actual service rate of the center. Unfortunately, the number of local requests in the system, *L*, will continue to increase until the system again stabilizes, so the local response times become unacceptable or the system runs out of some critical resource, for example the maximum number of connections.

⁴This assumes that all requests take μ_{act} seconds which is reasonable for a CPU or a disk drive.

MPL is a simple and effective control because it smoothes the external service demands by queuing requests outside of the provider, but MPL has several weaknesses. First, MPL is an integer value, so the impact on the local users cannot be fine tuned, and in particular, if an MPL of 1 causes too great of an impact, then external queries cannot be processed (see Sections 4.2 and 4.4). Second, MPL does not consider which resources will be affected, so the DBA must plan for the worst case of all the external requests accessing the same device (see Section 4.5). Third, small external requests get poor response times because they are forced to wait behind large requests (see Section 4.5). Finally, slow external consumers tie up provider connections and delay the processing of other external requests (see Section 4.6).

4.2 Fractional MPL via Spacing

If the response times of the local requests are still unacceptable even after an MPL is applied to the external requests, the requests can be further subdued by spacing job executions. By spacing jobs, the *average* MPL is reduced, thus allowing a fractional MPL to be specified. Spacing jobs further reduces the amount of time that local jobs must wait due to external jobs and gives the system time to recover from any short-term resource deficit caused by the external jobs, thereby avoiding long periods of lower QoS.

The interval between jobs can be specified in several ways. In its simplest form, a space can be a fixed time period. For example, after every job, pause for 10 seconds. A second way to define the space is based upon the amount of resources consumed by the job. The amount of resources used by the job can come from the provider, if the provider supplies such information, or the amount can be estimated by the profiler (Section 5.1). Defining spacing in this manner means that the pause after large jobs is longer than after short jobs. One last way to define spacing is based upon the amount of time the job spends in the provider. For example, if the space is specified as 100% and a job executes for 30 seconds, then the PA would wait an additional 30 seconds before allowing another job to enter the system. This last type of spacing not only responds to the length of the job, but it also responds to how busy the system is. Another advantage of this definition of spacing is that it correlates directly with MPL; the effective MPL is defined by

$$MPL_{eff} = \frac{MPL_{act}}{Spacing + 1}$$

For example, an MPL of 0.5 can be achieved with an MPL of 1 and a 100% spacing, and an MPL of 1.5 can be achieved with an MPL of 2 and a 33% spacing.

However, the effect of spacing is not identical to that of MPL, since the spacing oscillates more than MPL. While an external job is running, the local jobs experience higher response times, and while no job is running, the response times are lower. This oscillation implies that an MPL of 2 with 100% spacing does not control the system in the same way as an MPL of 1, but in either case, the average number of external requests running is one.

4.3 Suspending Jobs

Instead of spacing between external jobs, the jobs can be run more slowly by placing the spaces between sub-job units of work. The advantage of this technique is that a large job behaves like a group of small jobs, therefore the average local response time is achieved in a smaller time period. This technique also allows external jobs to use a processor sharing scheduling policy, although the quanta would be relatively large. We identified three ways of slowing down a job: intra-job spacing; suspending jobs; and dividing jobs. The effectiveness of these techniques depends upon the type of provider and the type of requests.

Intra-job spacing places an idle period between every block of the result. For example, the space could be placed between every 100 tuples fetched from a DBMS. If a result buffer exists between the provider and the PA, then intra-job spacing might not slow the external requests at all because the provider is filling the buffer while the job is sleeping. The buffering effect can be diminished by using larger block sizes. A more difficult problem with using intra-job spacing is the amount of data in the result might not correspond with the amount of work needed to find the result. For example, if an aggregate, like count, is applied to a large join query, the result is one tuple but the query might execute for an hour. The spacing should be based upon the amount of work done at the provider, rather than the size of the result.

Intra-job spacing can be generalized to job suspensions. When the system load is high, suspend the job until the system recovers. If a job is monopolizing an MPL unit for too long, suspend it and let another job run instead. Unfortunately, many providers offer no way to suspend a job except when its result buffer is filled, so jobs that cause constraint violations, for example, exceeding a maximum execution time, and cannot be satisfactorily suspended must be killed. Most providers can cancel a job without a long delay, although if the job is restarted, it will most likely need to be restarted from the beginning. A provider with a feature like the DB2 Governor Facility could assist the PA by automatically suspending the job after each unit of work.

Another technique that the PA can use to slow the execution of a job is to break the job into smaller jobs. The smaller jobs can then be scheduled using the job spacing of Section 4.2. Unfortunately, dividing jobs into smaller jobs is generally difficult or impossible.

The major drawback to all of these techniques is that they hold resources, like disk and memory buffers, locks, and connections, for a longer period of time. If a suspended job were holding read locks on some data and a local job wanted to update that data, then suspending the job actually slows down the local job, rather than speeding it up as intended. This problem still exists even when dividing the job into smaller jobs. If each of the smaller jobs are executed in a separate transaction, the combined answer could be inconsistent, so the jobs should be run within one transaction which implies that locks will be held between job executions.

Some providers, and in particular Oracle 7, have a form of optimistic concurrency control that allows updates to proceed in the presence of read locks. When a job tries to re-read a page that was updated since the last time it read it, the previous page is recovered from the log. If the job attempted to update the page instead of re-read it, the job would be aborted. This type of concurrency control allows read-only external jobs to be suspended without

any fear of blocking local jobs.

4.4 Feedback

In this section, we describe a feedback mechanism that monitors the utilization of key resources, for example, the disk drives, CPU, network, and memory, to decide when to start and stop jobs. This mechanism can be used in three related but distinct ways.

First, it can detect periods of relatively low activity, similar to the way the Condor system [18] finds idle resources to run batch programs. If the local workload experiences periods of high and low activity throughout the day, the PA can run external jobs during the periods of low activity. If an external job is running when the system is experiencing a heavy load, the PA can kill or suspend the job and restart it when the load decreases. For example, the PA could take advantage of low utilization during the lunch hour, a departmental meeting, or over the weekend.

Another way to take advantage of the utilization information is to detect the start and end of local jobs. If the local workload is such that periods of inactivity arise frequently and these periods are long enough for one or more external jobs to be completed, then the PA can start an external job whenever it detects one of these periods.

The third way that resource utilization can be used is to define the spacing of external jobs, similar to the spacing in Section 4.2. With this mechanism, external jobs are allowed to start only when the utilization of the busiest device is below some threshold. The mechanism can be extended so that if the utilization of some device exceeds another threshold, an external job is suspended or killed. The advantage of the feedback mechanism is that it responds more appropriately to changes in the local workload. For example, if the system is under a heavy load, feedback will keep the PA from starting a job, while spacing must occasionally start a job to determine that the system is indeed still busy.

The feedback mechanism is controlled by three parameters: the threshold that determines when to take action, the sampling interval, and the moving average time window. The utilization is sampled from the provider or the operating system by the PA once per sampling interval. Instead of basing decisions on one sample alone, a sliding window of samples can be averaged together. Using the moving average allows the PA to account for larger durations but the PA does not have to wait for the entire duration before making a decision.

The sampling interval controls the resolution of decisions within the PA. As the sampling interval gets smaller, the PA gets a more detailed view of the system. If the sampling interval is set too low, the overhead of sampling will adversely affect performance. Conversely, if the sampling interval is too high, the PA becomes sluggish because as far as the PA is concerned, the state of the provider has not changed. The proper setting for the sampling interval depends upon the duration of the external jobs and the volatility of the local workload.

The moving average time window controls how responsive the PA is to changes in the provider. A small time window means the PA will react quickly, but if the window is too small the PA tends to over-react. If the time window is longer, the PA is more tentative and responds slowly to both increases and decreases in utilization. When using feedback to detect changes in the workload, the moving average time window should be long enough to get a

reasonable estimate of the true average and filter out the variation. If the goal is to detect when no local jobs are using the system, the window should be small, perhaps eliminating the moving average altogether. When controlling the job spacing, the size of the window depends upon the duration of the external jobs. The window should be long enough to at least cover the execution of the job and the desired spacing. When using feedback to stop external jobs, the window should no less than the maximum allowed duration of an external job.

4.5 Profiling

Information about the potential resource consumption of an external job is invaluable to the PA, but when an external requests arrives it is tagged with any such information. As discussed in Section 4.1, one of the major drawbacks of limiting the MPL is that small jobs are forced to wait behind large jobs. But by using the job profile information, the PA can schedule small jobs differently from large jobs. From the perspective of protecting the local users, scheduling small and large jobs differently is intuitively reasonable because a single small job is not as likely to seriously impact the performance of the system.

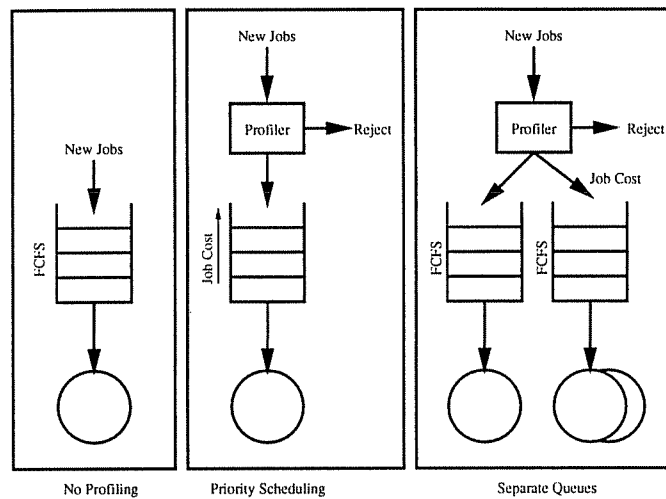


Figure 2: Queuing Strategies

The profile information can be used in a variety of ways (see Figure 2). One possible strategy groups jobs that perform less than 10 I/Os into one group, between 10 and 100 into a second group, between 100 and 10,000 in a third group, and rejects any job that is expected to do more than 10,000 I/Os. Each group is given its own queue and scheduling policy (MPL, spacing, etc.). The maximum MPL in this case is the sum of the MPLs for each group.

Another strategy uses a single queue and scheduling policy, but orders the queue by the number of expected I/Os⁵ With this strategy, a small job would wait behind at most one job that was longer than it, and only because the longer job was already executing when the small job arrived. If the MPL were more than one the small job must wait for only one job to complete.

⁵This strategy is often called shortest-job first or shortest processing time first in the scheduling literature.

The advantage of the first strategy is that small jobs will get better service because they never wait for a large job, but the second strategy is a little easier for the administrator to configure. Also, if the administrator decides that only one external job may be executing at any given time, only the second strategy is feasible.

Notice that with both of the strategies, absolute accuracy of the profile is not necessary. As long as the profiles are accurate relative to one another, both schemes function properly, except when a job is rejected. Occasionally, a profile could be completely wrong. The PA must ensure that the use of an incorrect profile will not detrimentally affect the protection that it provides, and secondarily, the PA should keep small jobs from waiting behind a large job that masqueraded as a small job. One way to provide insurance against running a job with a bad profile is to place a time limit on the execution of a job where the time limit is based upon the job's profile. The PA should also provide override that allows the administrator to accept a request that was mistakenly rejected.

When the profile contains details regarding the individual resources that a job is expected to use, the PA can more fully utilize the the providers resources. Either of the above strategies could be extended, for example, with a policy that allowed up to 5 jobs to run concurrently (MPL=5) as long no two jobs accessed the same disk drive. Whenever less than 5 jobs are running, the PA looks through its queue and starts the first job that doesn't access any of the devices that are already in use.

Most data providers quickly estimate the cost of executing a job without actually executing it. For example, most relational DBMSs use an optimizer to estimate the cost for candidate execution plans with the goal of finding the least expensive plan. The optimizers usually compute the cost as a weighted sum of the expected number of CPU cycles, disk I/Os, and network messages. Using additional catalog information about the location of relations, the total disk I/Os can be broken down on a per disk basis. Although optimizers do not need this information, the PA could use it to keep from over utilizing any one device.

Unfortunately, the profile information is usually not exported by the provider — most DBMSs discard the cost estimates once the plan is produced. The PA could use the plan and reapply the statistics that the DBMS used to produce the plan to reconstruct the profile, but often the statistics are not exported either⁶.

A work-around to these limitations is to replicate the profiler in the PA, similar to the approach taken in the Pegasus project at HP [10]. The profiler will frequently consult catalog information, so to avoid placing an additional strain on the provider, all of the needed catalog information should also be replicated within the PA. The advantage of using the provider's profile is that no effort is duplicated and the profile will generally be more up-to-date than the PA's profile, for example when an index is added or dropped but the PA has not been updated. The problem with using the provider to produce the profile, when it is an option, is that requesting profiles from the provider consumes provider resources. In some cases, for example, a 15-way join, the amount of resources used is far from trivial, implying that profiling must also be scheduled.

⁶This is especially true when the DBMS maintains histograms rather than simple averages.

4.6 Buffering

Until now, we have not considered the rate at which the external users consume their results. We implicitly assumed that the results were consumed as quickly as they were produced. What happens when results are consumed slowly? Assume that the only control was an MPL of 1. If one job decided not to consume any of its results for a while, the PA could not process any additional jobs and the job would be holding resources, like locks, for a long time. Similarly, since the PA is can be used on the Internet, the network connection between the PA and the external user could be slow. Again, the other external requests will pay for the added delays caused by this request.

One way to limit the impact of a slow consumer is to enforce a maximum execution time or a minimum transfer rate. If the constraints are violated, the job is suspended or killed and another job takes its place. A more forgiving strategy places a buffer between the PA and the external client. If the buffer is large enough, the PA can pull the entire result out of the provider as quickly as possible and allow the job to proceed as slowly as it wishes. The provider is no longer impacted by the slow consumer, and the PA can execute another job. Of course, the same problem could happen with the next job, and eventually the PA would run out of buffer space. When faced with this situation, the PA must either kill one of the slow jobs and reclaim its buffer space, or the PA must stop executing jobs and wait for enough buffer space to clear.

The buffer between the PA and the client can be both main memory and disk pages. All of the buffer space can be pooled together and shared among all of the PA clients. The page size should be relatively large since all I/O within one connection will be sequential. When swapping out a page to disk within one connection; the most recently filled page should be sent to disk. When choosing a page to swap out from all connections, the page least likely to be accessed should be swapped out. To find the least likely page, let r_i denote the rate at which connection i is consuming pages, and d_i denote the distance that the last page of a connection is from the head of the queue. Then $t_i = d_i/r_i$ denotes the expected time until the last page is referenced. The least likely page to be accessed is the page with the maximum expected reference time.

4.7 External Priorities

The PA can be extended to support priorities among external requests. The techniques used to enforce priorities are similar to the techniques used in Section 4.5 to take advantage of profile information. One queue ordered by priority can be used to essentially achieve absolute priorities. Multiple queues, one per priority, with the total MPL divided among the queues based upon priority achieves relative priorities. The other mechanisms can also be extended with priorities. For example, when a job must be suspended or killed, choose the job with the lowest priority.

4.8 Provider Specific

Each provider comes with its own set of features that could be used to assist the PA in its task. For example, some providers have separate buffer pools based upon user-id, some have prefetching, and others have priority scheduling.

While we encourage the use of these types of features, this paper focuses on general solutions, so these features are not discussed further in this paper.

5 Implementation

We implemented a prototype PA to experiment with the mechanisms and evaluate policies described in the previous section. This section describes the design decisions we made while implementing the PA for a modern commercial relational DBMS.

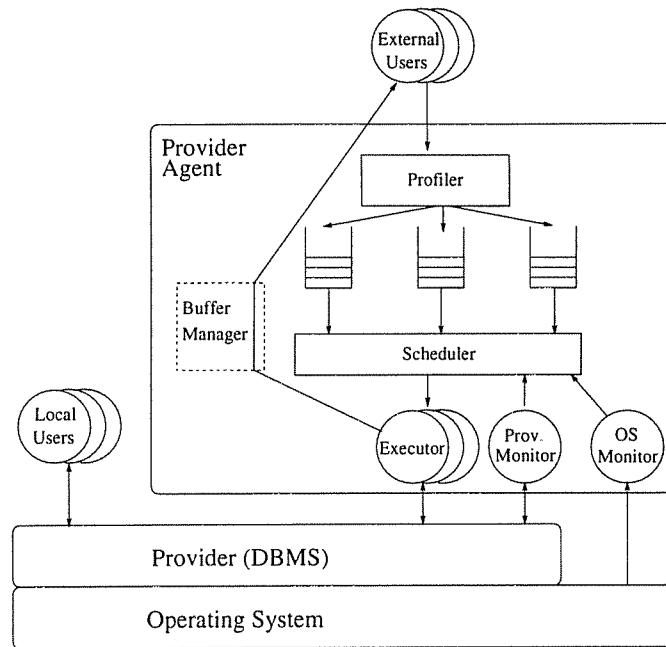


Figure 3: Provider Agent Architecture

5.1 Profiler

When a request arrives at the PA, it is first profiled to estimate its resource needs and the expected impact that this request will have on the provider. The profiler is a fully functional query optimizer designed to mimic the provider’s optimizer. We chose to replicate the query optimizer inside the PA because the DBMS that we considered, like most, did not export the cost estimate. The goal of the profiler is not to find the best plan, but the plan that is expected to be produced by the provider. If the PA maintains additional statistics that the provider does not consider, the statistics should not be used during plan generation, but once the plan is chosen the PA can use the additional statistics to produce a better profile.

The profiler will frequently consult catalog information, so to avoid placing an additional strain on the provider, we cached all of the catalog information within the PA. Although not implemented in our prototype, the PA must

have some mechanism for updating its statistics and becoming aware of new database objects. The mechanism could be automatic updates, or the administrator could update the catalogs as needed.

5.2 Scheduler

After a request is profiled, it is queued until the scheduler decides that the request can be safely executed without adversely affecting the local requests. This component uses the mechanisms described in Section 4 to enforce the policies established by the administrator, and in the next section, we evaluate the effectiveness of several possible schedulers.

5.3 Executor

An executor is the provider client that processes requests on behalf of the PA and its external clients. The executor contacts the external client to inform it that its job is ready to be processed. The executor then sends the request to the provider and receives the result which it translates to a standard form sends it to the client, either directly or through the buffer manager. Our current prototype does not include a buffer manager, and in our experiments, external clients consumed results as quickly as possible.

5.4 System Monitors

The monitors periodically sample statistics provided by the operating system and the provider to inform the scheduler about current resource utilizations. The OS monitor was the only component of the PA that ran on the same machine as the provider because the operating system we used did not allow remote applications to directly obtain OS statistics.

6 Experiments

In this section, we evaluate the effectiveness of the key mechanisms of the prototype PA. We demonstrate the need for local QoS protection by running a moderately heavy, but realistic external workload without the PA. Although we can disrupt local processing to an arbitrary degree, we show that even a moderately heavy load can increase average response times over 2,000%.

We then add the PA and show that load control with spacing allows the DBA to set the worst case performance of the local workload. Once we establish that the basic load control mechanism is effective at controlling the performance degradation, we improve the PA by using a dynamic load control mechanism. We show that the dynamic control responds to the local workload by allowing more external jobs to be processed when the system is under-utilized. Lastly, we demonstrate that profiling dramatically improves the throughput of small external jobs in a mixed external workload.

We generated a synthetic workload for the experiments so that our results were repeatable, but we used our prototype PA and a current version of a commercial DBMS to ensure that our results were realistic. The local workload modeled a transaction processing system. We ran one experiment that performed updates to verify that updates did not affect our results unduly, but in the rest of the experiments the transactions were read-only. The external workload modeled exploration queries that needed to perform scans over large portions of the data. The experiment parameters are summarized in Table 1.

Table 1: Experiment Parameters

Class	Parameter	Setting
local	#terminals	40
	inter-arrivals	exp(0.2) sec.
	queries	10 compiled point selects, using non-clustered index
external	MPL	0 - 20
	#jobs	infinite
	queries	range count, 25-75% of relation, using clustered index
system	#disks	1
	#relations	7
	tuple size	208 bytes
	relation size	12MB
	data cache	7.4MB of 1 page blocks 5.0MB of 8 page blocks
	duration	25 minutes + 5 minutes warm-up without external

6.1 Experiment 1: No Control

To demonstrate the impact that external requests can have on the local response times, we ran a constant external workload of 20 queries. Figure 4 shows the cumulative distributions of the local response times over the 25 minute experiment. An MPL of 20 increased in average response over 20 times more than the response when no external jobs were run. Although not illustrated by an experiment, when we pushed the system harder, we inadvertently caused a large percentage of local requests to fail because of a lack of memory needed to run the queries. Notice even a small period of high external activity is enough to wreak havoc on the hapless local users.

6.2 Experiment 2: Fixed MPL

In Figure 5, we show that the local response time can be controlled by adjusting, via MPL and spacing, the average number of external requests in the system. Decreasing the MPL shifts the entire distribution of response times

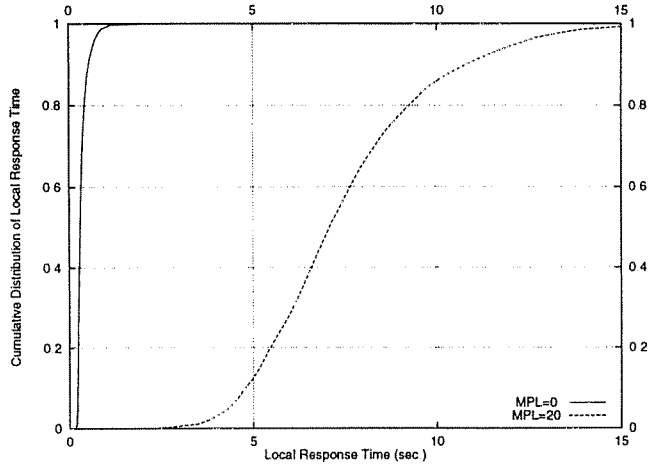


Figure 4: Impact of heavy external load

towards the original response time curve when no external queries were run (MPL=0). For this particular external workload, increasing the MPL beyond 1 actually decreases the throughput of the external requests because the competition between external jobs changes the disk access pattern from sequential to random access.

We ran a similar experiment except that each local transaction updated the last record read after a 3 second delay. Figure 6 shows that the external MPL still controls the local response time, therefore the techniques described still function properly in the presence of updates.

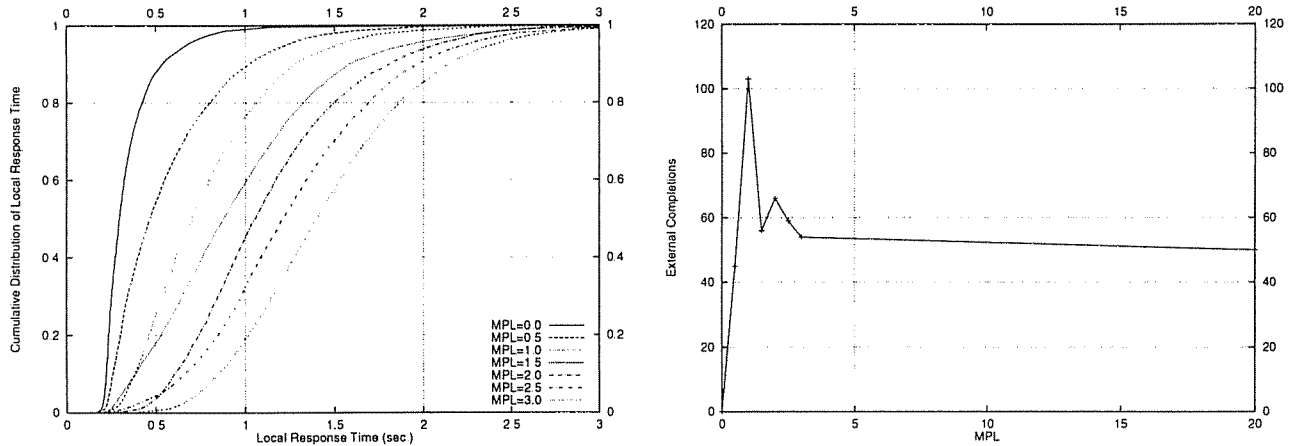


Figure 5: Fixed MPL Control

6.3 Experiment 3: Dynamic MPL

Figure 7 illustrates that by adjusting the add threshold, monitoring the utilization to control the spacing of external requests protects the local response times much like spacing. The advantage of this feedback mechanism is that it can aptly avoid periods when the local users place a heavy load on the system, or take advantage of periods of

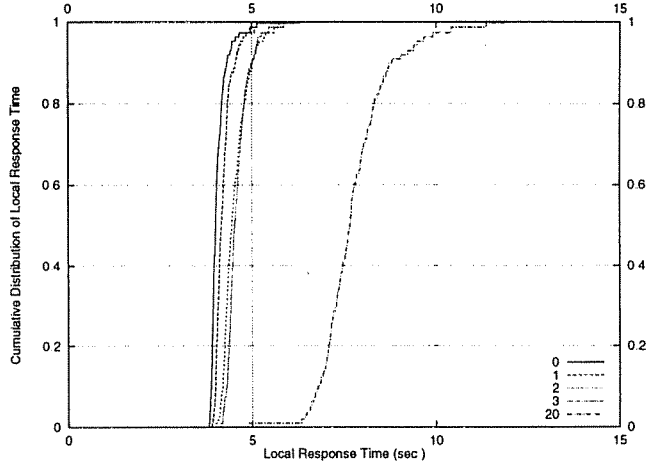


Figure 6: Fixed MPL with Updates

low activity. Figure 8 shows the difference between spacing and feedback. When we increased the average local interarrival time from 0.2 seconds to 0.4 seconds (i.e., decreased the local demand), the dynamic version kept the local users at the same QoS, but allowed more external work to be completed than spacing.

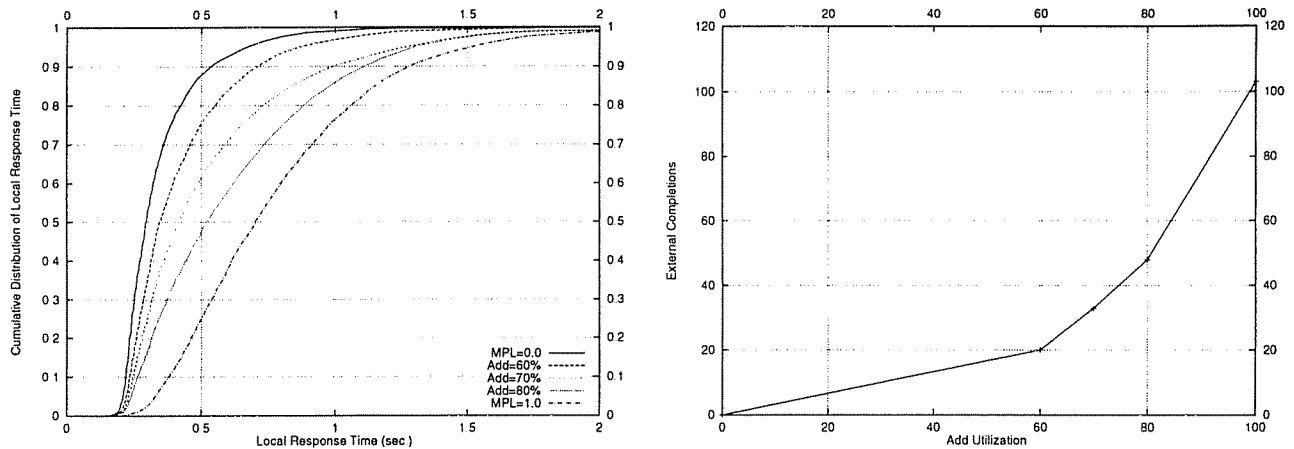


Figure 7: Dynamic MPL Control

6.4 Experiment 4: Profiling

Profiling allows the PA to keep small jobs from waiting for large jobs to complete. This experiment, which is summarized in Table 2, illustrates the need for profiling by running an external workload with both large and small jobs. The arrival times of the queries were uniformly distributed over the 25 minute experiment. We assume that the administrator decided that they wanted at most one large query in the system at a time. Furthermore, they chose a dynamic MPL with an add utilization of 70%. In the first run, the profiler was not used, so small queries entered the same queue as the large queries. In the next run, the profiler was used, and the administrator decided

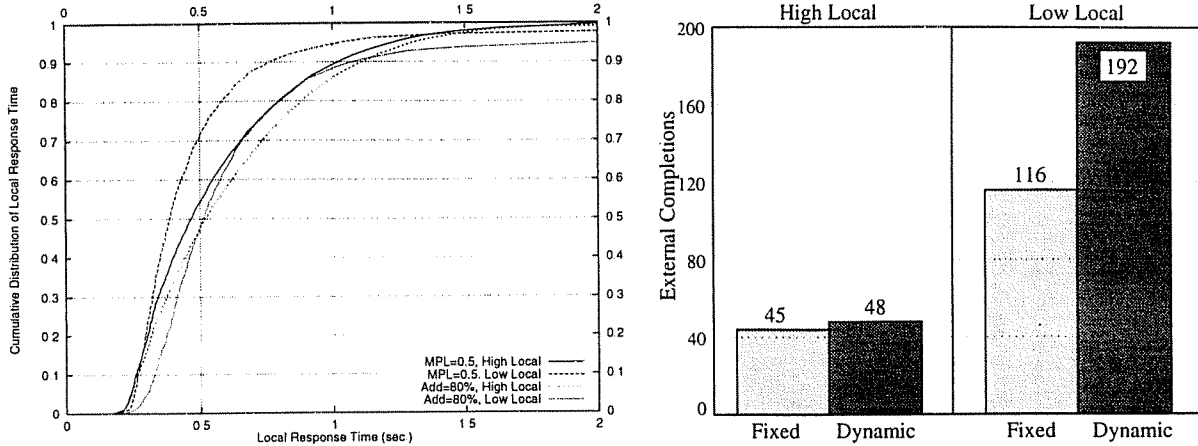


Figure 8: Dynamic vs. Fixed MPL

that at most two small queries could enter the system at a time as long as the utilization was below 90%. So in this case, the PA scheduler consisted of two queues, one for small queries and one for large queries.

The results of the experiment are shown in Figure 9. The left chart shows that allowing the two small queries to execute along with the large query had little impact on the local QoS. The chart on the right, however, shows that the throughput of the small queries increased by 25 times when profiling was used, with only a small decrease in the throughput of the large queries.

Table 2: Mixed External Workload Parameters

Class	Parameter	Setting
large external	MPL	1
	add utilization	70%
	#jobs	76
	queries	range count, 25-75% of relation, using clustered index
small external	MPL	2
	add utilization	90%
	#jobs	2530
	queries	point select, using non-clustered index

7 Related Work

Our work inherits much from the work on operating system load control [14, 9]. Our contribution is in applying load control in a novel way to protect the local QoS. A significant amount of work on database scheduling has also been completed, especially on memory allocation [15, 13, 19, 20, 3, 4, 8]. The main difference between our work and the

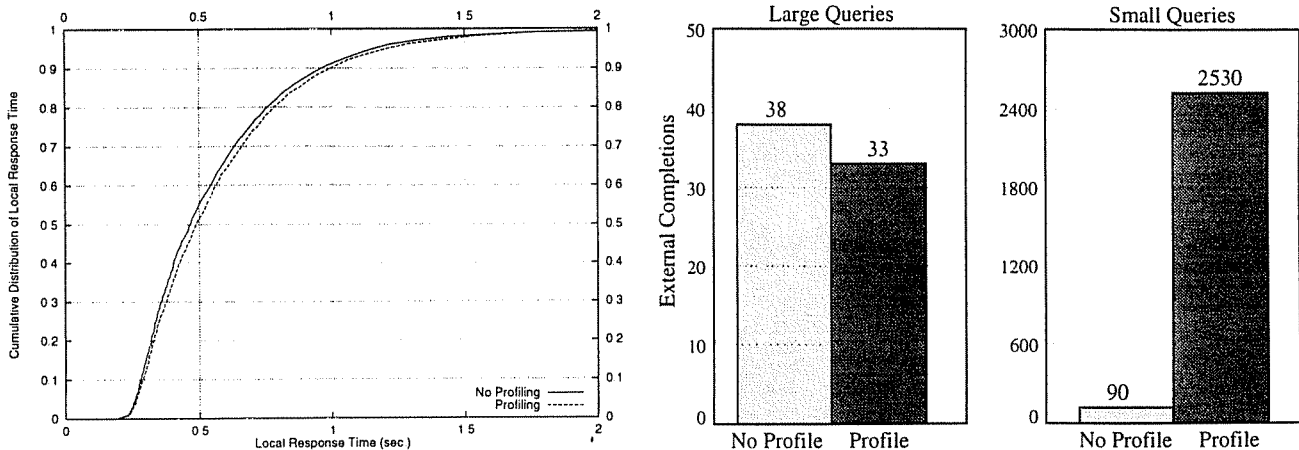


Figure 9: Profiling Results

database scheduling is that ours exists outside the database and is not specific to any particular DBMS.

A few vendors, especially noteworthy is IBM's DB2 for MVS [7], have some load control mechanisms similar to ours, but as far as we are aware, they do not have our fractional MPL mechanisms and they do not use profiling to enhance external scheduling. Two other classes of applications that extend the scheduling of DBMSs are currently available: Query Analyzers and Transaction Processing (TP) Monitors.

Query analyzer products are sold, for example, by Platinum Technologies (Query Analyzer, DB Analyzer, SQL-Spy, Governor Facility, Detector) [23], Blue Lagoon (DBProfiler), Level\5 (SmartMode, Quest) [17], and MicroStrategy (DSS Administrator) [21]. The analyzers provide a feature like our profiler to estimate the cost of executing a query before actually running it. The difference between query analyzers and our project is that the analyzers use this information to warn users and developers of poorly formed queries that could consume large amounts of resources, while we, on the other hand, use the profile for scheduling. An exception to this is DSS Administrator which appears to use the profile information to schedule the queries from their Decision Support System (DSS) product.

TP Monitors are sold by Transarc (Encina [12]), Novell (Tuxedo [1]), NCR (Top End [25]), and IBM (CICS [6]), to name a few. The main thrust of a TP monitor is to coordinate transactions through a TP system, but TP monitors do offer priority queuing and load balancing [2]. The priorities controlled by the applications, not through profiling, and their load balancing just tries to keep the machines in the TP system equally busy. NCR is developing a product called DBQM that addresses the issue of external query access, but we were unable to get additional information before submission.

8 Future Work

We have identified a number of areas for future work. We described some mechanisms in Section 4, for example the buffer manager, that we believe to be useful, but we have not yet demonstrated that fact with experiments. Another issue with our current prototype is that it can be unfair to large requests because small requests can keep

large requests from starting. The grouping of multiple requests into a transaction, from the perspective of locking and resource consumption, presents additional challenges to the PA. We are also looking into the ParaDyn project [22] to see how it can enhance the performance monitoring of the PA.

9 Conclusions

We identified the potential performance dangers of allowing external access to an existing information system, and we described several such scenarios that need a solution. We believe organizations should address this problem before allowing external access. We offer the PA as a general solution that can, with little effort, protect most sites without requiring changes to the underlying system or any of the programs that the site uses for local access.

We described how load control, request profiling, buffering, and priority scheduling can be combined to form an elegant, novel solution. We demonstrated the need for load control, and that MPL and spacing effectively limit the impact of external requests, even in the presence of updates. Feedback is used to obtain a dynamic MPL that allows the PA to respond to changes in local workload. Profiling enables the administrator to base scheduling decisions on the resource requirements of a job. In particular, we showed that discriminating between large and small jobs improved the processing of the small requests by over 25 times while only marginally slowing the large requests. All these features working from the outside of the database combine to provide a realistic and effective solution to a very real problem.

We believe that the demand for the PA will continue to increase as Internet services, data mining, and data integration projects proliferate. We also believe that the ideas presented in this paper can be quickly integrated with existing products like middleware tools, TP monitors, query analyzers, multi-database systems, and database schedulers which implies that commercial products with PA-like features should appear in the near future.

References

- [1] Juan M. Andrade. Open on-line transaction processing with the TUXEDO system. *IEEE CompCon*, 1992. also <http://www.beasys.com/Product/tuxedo.htm>.
- [2] Philip Bernstein. Transaction processing monitors. *Communications of the ACM*, 33(11), November 1990.
- [3] Kurt P. Brown, Michael J. Carey, and Miron Livny. Managing memory to meet multiclass workload response time goals. In *Proceedings of the 19th VLDB Conference*, Dublin, Ireland, 1993.
- [4] Kurt P. Brown et al. Towards automated performance tuning for complex workloads. In *Proceedings of the 20th VLDB Conference*, Santiago, Chile, 1994.
- [5] Sudarshan Chawathe et al. The tsimmis project: Integration of heterogeneous information sources. In *Proceedings of IPSJ Conference*, pages 7–18, Tokyo, Japan, October 1994.

- [6] IBM Corp. CICS product family home page. <http://www.hursley.ibm.com/cics>.
- [7] IBM Corp. MVS workload manager / system resource manager.
<http://www.s390.hosting.ibm.com/products/mvs/wlm/index.html>.
- [8] Diane L. Davidson and Goetz Graefe. Dynamic resource brokering for multi-user query execution. In *SIGMOD 95*, pages 281–292, San Jose, CA, 1995. ACM.
- [9] P.J. Denning. The working set model for program behavior. *Comm. ACM*, 11(5):323–333, May 1968.
- [10] Weimin Du, Ravi Krishnamurthy, and Ming-Chien Shan. Query optimization in heterogeneous dbms. In *Proceedings of the 18th VLDB Conference*, pages 277–291, 1992.
- [11] David Eichman. Ethical web agents. In *Electronic Proceedings of the Second World Wide Web Conference '94: Mosaic and the Web*, 1994.
<http://www.ncsa.uiuc.edu/SDG/IT94/Proceedings/Agents/eichmann.ethical/eichmann.html>.
- [12] Jeffery L. Eppinger and Scott Dietzen. Encina: Modular transaction processing. *IEEE CompCon*, 1992. also
<http://www.transarc.com/afs/transarc.com/public/www/Public/ProdServ/Product/Encina/index.html>.
- [13] C. Faloutsos, R. Ng, and T. Sellis. Predictive load control for flexible buffer allocation. In *Proceedings of the 17th Int'l VLDB Conf.*, Barcelona, Spain, September 1991.
- [14] Domenico Ferrari. *Computer Systems Performance Evaluation*. Prentice Hall, 1978.
- [15] Robert Brian Hagmann and Domenico Ferrari. Performance analysis of several back-end database architectures. *ACM Transactions on Database Systems*, 11(1):1–26, March 1986.
- [16] Martijn Koster. Guidelines for robot writers. <http://info.webcrawler.com/mak/projects/robots/guidelines.html>.
- [17] Level\5. <http://www.l5r.com/>.
- [18] M. Litzkow, M. Livny, and M. W. Mutka. Condor - a hunter of idle workstations. In *Proceedings of the 8th International Conference of Distributed Computing Systems*, pages 104–111, June 1988.
- [19] Manish Mehta and David J. DeWitt. Dynamic memory allocation for multiple-query workloads. In *Proceedings of the 19th VLDB Conference*, Dublin, Ireland, 1993.
- [20] Manish Mehta, Valery Soloviev, and David J. DeWitt. Batch scheduling in parallel database systems. In *Proceedings of the 9th International Conference on Data Engineering*, Vienna, Austria, April 1993.
- [21] MicroStrategy. <http://www.strategy.com/>.
- [22] Barton P. Miller et al. The paradyn parallel performance measurement tools. *IEEE Computer*, 28(11), November 1995. Special issue on performance evaluation tools for parallel and distributed computer systems.

- [23] Platinum Technologies. <http://www.platinum.com/>.
- [24] Avi Silberschatz, Mike Stonebraker, and Jeff Ullman. Database research: Achievements and opportunities into the 21st century. *SIGMOD Record*, 25(1):52–63, March 1996.
- [25] Randy Smerik. An overview of TOP END. *IEEE CompCon*, 1992. also <http://www.ncr.com/product/topend>.
- [26] Michael Stonebraker et al. Mariposa: a wide-area distributed system. *VLDB Journal*, pages 48–63, 1996.