

Self-tuning Page Cleaning in DB2

Wenguang Wang Richard B. Bunt
Department of Computer Science
University of Saskatchewan
57 Campus Drive
Saskatoon, SK S7N 5A9
Canada
Email: {wew036, bunt}@cs.usask.ca

Abstract

Storage management is an important part of DB2. The buffer pool in DB2 is used to catch the disk pages of the database, and its management algorithm can significantly affect performance. Because of the complexity of DB2 and the workloads running on it, the buffer pool management algorithm is hard to study, config, and tune. In order to investigate the buffer pool management algorithm under controlled circumstances, a trace of buffer pool requests was collected and a trace-driven simulator was developed. The impact of various parameters of the buffer pool management algorithm was studied in the simulator. Relationships among different activities competing for storage space and the I/O channel were examined. A self-tuning algorithm for buffer pool management was developed and tested in both the simulator and DB2. Simulation tests showed that the new algorithm can achieve comparable performance to a hand-tuned system. The experiments in DB2 on a small and a medium TPC-C database verified the simulation results.

1 Introduction

In a Database Management System (DBMS), a buffer pool is used to cache the disk pages of

the database. Because the speed gap between disk and memory is large, the effectiveness of the buffer pool management algorithm is very important to the performance of the DBMS. Tuning such an algorithm is a complex task because of the number of parameters involved and the complexity of their interactions. It requires a detailed understanding of the nature of activities that compete for the resources being managed, including storage space and the I/O channel. Because of the complexity of a commercial DBMS, however, it is often difficult to analyze the buffer pool algorithm, or implement a new one and test it directly. Simulation provides an effective alternative.

For our research, a blend of direct experimentation and trace-driven simulation is being used in a detailed study of buffer pool management in DB2, IBM's popular commercial relational DBMS. The TPC-C [9] benchmark provides the workload. A detailed buffer pool simulator has been developed for DB2 and verified. To provide a realistic reproducible workload for the simulation experiments, a trace of the buffer pool requests when running the TPC-C benchmark was captured. The effects of parameters of the buffer pool algorithm are examined by simulation. Activities of the I/O channel in the buffer pool are investigated as well. A self-tuning algorithm was proposed to simplify the configuration of the number of

page cleaners of the buffer pool management algorithm. Simulation results and experiments in DB2 showed that the new algorithm can achieve comparable performance to hand-tuned system.

This paper is organized as follows: §2 gives a brief review of some previous work relating to database buffer pool management and self-tuning architecture; §3 describes the buffer pool management algorithm used in DB2; §4 presents the trace collection method and the trace formats; §5 discusses the simulator and its verification; §6 presents the results of the experiments conducted to study the relationships among buffer pool I/O activities and their impacts on the buffer pool performance; §7 discusses a self-tuning algorithm and its experiment results in simulation and measurements; §8 shows the simulation results and the real system experiment results of the self-tuning algorithm; §9 presents conclusions and future work.

2 Related Work

It has been observed [7, 5, 8] that the references to the disk pages of a database have locality. Furthermore, this locality of reference is more regular, predictable, and exploitable than the localized reference activity found in programs in general. For this reason, a buffer pool is employed by a DBMS as a cache of disk pages. Various database buffer pool management algorithms, including LRU, FIFO, CLOCK, LRD, DGCLOCK, and Working Set [3], are analyzed and compared in [4]. A set of flexible predictive buffer management algorithms was proposed in [6].

All approaches aim to improve the hit ratio of the buffer pool in order to reduce disk reads. However, another problem is also important—how and when to clean out dirty pages (i.e. writing back to disk pages that have been changed since being fetched into the pool). Two approaches can be applied: page cleaning on demand and asynchronous page cleaning. In page cleaning on demand, a dirty page will be cleaned only when it is selected for replacement. In asynchronous page cleaning, dirty pages are cleaned asynchronously before

they are selected for replacement. In traditional virtual memory systems, processes write only to their data segments, and normally not many dirty pages are generated. However, in a database buffer pool, transaction-based workloads may generate so many dirty pages that page cleaning on demand is not efficient. The effect of the asynchronous page cleaning strategy is investigated in this paper, and the approach proposed can be applied to all algorithms mentioned above.

Because of the complexity of DBMS tuning, some goal-oriented tuning algorithm are proposed. [2] presents an approach to dynamically adjust the buffer pool sizes of DBMS based on response time goals. [1] proposes a goal-oriented tuning architecture to convert the low level control knobs (buffer pool size, working buffer size, etc.) to high level transaction response time goals.

3 The DB2 Buffer Pool Management Algorithm

The buffer pool management algorithm contains several parts: fetching, placement, and replacement. The fetching algorithm brings new pages into the buffer pool. Normally, fetching on demand is used, but when there are sequential accesses to pages, prefetching can be a helpful supplement to fetching on demand. DB2 supports both fetching on demand and prefetching. For this analysis, only fetching on demand is used because there is no sequential access in the TPC-C workload. The placement algorithm defines how to place a page into the free pool of the buffer. The replacement algorithm defines how to maintain the free pool of the buffer. Replacement on demand makes a free space only when needed. Pre-replacement tries to always keep spaces available by writing back to disk dirty pages that are deemed not to be needed. Traditional buffer management algorithms use only replacement on demand, but both approaches are used in DB2 buffer pool management. By performing page cleaning, the buffer pool maintains some space holding clean pages that can be used immediately for new pages.

3.1 Basic Replacement Algorithm

Unlike traditional virtual memory management, there is no hardware support in the buffer pool to access a page that is not in memory. Therefore, a *fix/unfix* mechanism is used. When the DBMS needs to access a page, it will send a *fix* request to the buffer pool. If the page is already in the buffer pool, no physical I/O is needed; otherwise, it is read into the buffer pool from disk. The DBMS can access this page freely from the buffer pool and it cannot be evicted from memory after the *fix*. When the DBMS finishes using this page, an *unfix* request is sent to the buffer pool. After the *unfix*, this page is allowed to be evicted when needed.

3.2 Page Cleaning Algorithm

When a page is selected for replacement, its status is checked. If it is clean (unchanged), its space can be used immediately. If it is dirty, a *synchronous* write must be performed to clean this page. Synchronous writes negatively impact the throughput of the system because the space occupied by this dirty page cannot be used until the synchronous write completes.

To alleviate this problem, one or more *page cleaners* are used to clean out pages asynchronously before synchronous writes are needed. At the beginning, all page cleaners are sleeping. Three kinds of events can wake up the page cleaners:

- *Exceeding the Dirty Page Threshold.* The dirty page threshold is a configurable parameter which indicates the desired percentage of dirty pages in the buffer pool. Page cleaners will be awakened when the percentage of dirty pages exceeds the threshold. The default value of the threshold is 60%.
- *Dirty Replacement.* When a dirty page is selected from the buffer pool for replacement, i.e. a synchronous write occurs, the page cleaners wake up.
- *Exceeding the softmax value.* When the percentage of changes recorded in the log

file exceeds the softmax value, the page cleaners wake up.

When a page cleaner wakes up, it collects some dirty pages and writes them back to the disks. If the condition for page cleaning still exists, the page cleaner will begin another round of collecting and cleaning. If the conditions for page cleaning do not exist, the page cleaner will sleep. When a page cleaner writes pages, the clients do not need to wait for the writes directly. Therefore, these writes are called asynchronous writes.

Figure 1 shows the I/O activities between the buffer pool and the disks. The buffer pool can be considered as two regions: the *clean region* which contains clean pages, and the *dirty region* which contains dirty pages. When the applications are running, many pages are read from the disks. At the same time, many pages in the buffer pool are changed, which causes the dirty region to expand. The synchronous writes and asynchronous writes write the dirty pages back to the disks which causes the dirty region to shrink. When the system is in stable state, the number of pages changed is equal to the number of pages cleaned so that a constant supply of free pages is available.

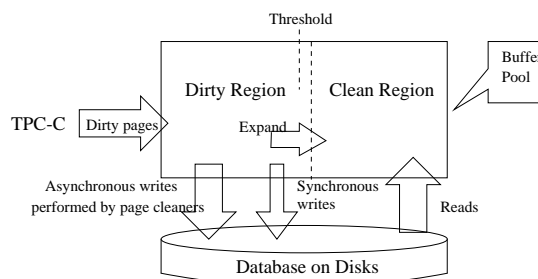


Figure 1: Buffer Pool and Its I/O Activities with Disks

4 Trace Collection

4.1 Workload – TPC-C Benchmark

The TPC-C benchmark is used to provide the workload to the simulator and to provide the “real” system against which simulation results

are compared. This benchmark is among a group of benchmarks defined by the Transaction Processing performance Council¹ (TPC). The TPC-C Benchmark is an On-Line Transaction Processing (OLTP) workload. It is a mixture of read-only and update-intensive transactions that simulate the activities found in complex OLTP application environments. Multiple transactions are used to simulate the business activity of processing an order, and each transaction is subject to a response time constraint.

There are five kinds of transactions in the TPC-C benchmark, and these are listed in Table 1. The percentage column shows their respective percentages in the total number of transactions. For example, 45% transactions are New-Order transactions. Because the New-Order transaction is the backbone of the workload, the performance of the TPC-C benchmark is expressed in Transactions Per Minute (TPM), which is defined as the number of New-Order transactions processed per minute. The size of the TPC-C database is given by the number of “warehouses” defined. There are about 100M bytes of data for one warehouse.

Table 1: TPC-C Transactions

Name	Percentage
New-Order	45%
Payment	43%
Order-Status	4%
Delivery	4%
Stock-Level	4%

4.2 Trace Collection and Contents

In order to collect the trace of the buffer pool activities needed by the simulator, a TPC-C database with 50 warehouses was created. The benchmark runs on a PC Server running Windows NT Server 4.0 in the Distributed Systems Performance Laboratory at the University of Saskatchewan. The DB2 version used is 6.1. The TPC-C database is installed on 9 disks.

When the TPC-C benchmark is running, transactions are sent to the DBMS. The DBMS

¹<http://www.tpc.org/>

executes the transactions and requests disk pages from the buffer pool. All fix/unfix requests to the buffer pool are needed for the simulation. The information of every request going through the trace point shown in Figure 2 is recorded in a trace buffer in memory by the trace tool. Because the trace buffer cannot hold all the trace information needed for the simulator at one time, the TPC-C program was changed so that it suspends before the trace buffer is full, dumps the trace into a file, and continues to run. By using this technique, arbitrarily large traces can be obtained.

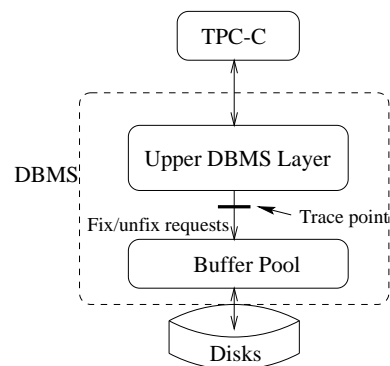


Figure 2: Trace collection structure

Because of the overhead associated with this data collection, the system runs more slowly when recording the trace. Therefore, the sequences of the aggregated trace records from all agents may change, but the sequence for each agent does not change. The trace used by the simulator is separated based on the agent ID. Therefore, although the aggregated trace sequence changes, it will not affect the simulation result. The begin and end time of each transaction is recorded in the TPC-C program. This transaction information is inserted into the trace recorded from the buffer pool sorted by the timestamps. After these processes, the trace can be used by the simulator. The trace used in these experiments includes about 60 million buffer pool requests from 200 thousand transactions.

Necessary information about fix and unfix requests is recorded in the trace. Table 2 shows the important fields of a trace record. The *fix mode* indicates the type of fix request. There are two kinds of fix: *exclusive* and *shared*. Any

Table 2: Important Fields of the Trace

Field	Value
type	The type (fix or unfix) of this record
tid	The number to identify requests from different clients
object type	Can be INDEX or DATA
object id	The number to identify the table this page belongs to
page number	The disk address of the page
fix mode (only for fix)	Can be EXCLUSIVE or SHARE
modify flag (only for unfix)	Indicates whether or not this page has been modified

page can have at most one exclusive fix but can have multiple shared fixes at the same time.

5 The Buffer Pool Simulator

The buffer pool simulator simulates DB2's buffer pool management algorithm and the disk subsystem. This simulator contains about 8000 lines of C++ code. An event-driven architecture is used in the simulator. Different components of the simulator communicate through events. Figure 3 shows the components and main event types.

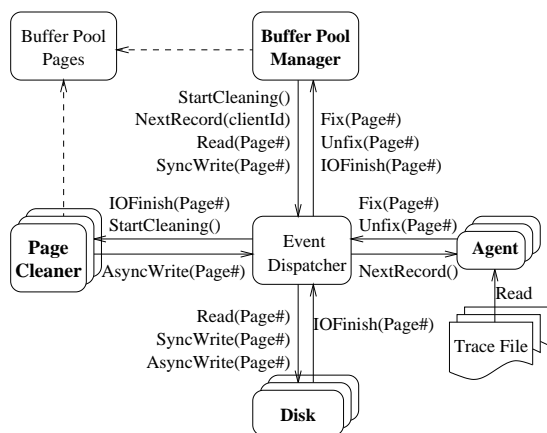


Figure 3: The Structure of the Simulator

There are four basic components in the simulator.

- The *Buffer Pool Manager* contains the basic buffer pool management algorithm. It manages the placement and replacement

of the buffer pool pages. It accepts the fix/unfix events, and sends out disk I/O events (Read and SyncWrite). It also sends StartCleaning events to trigger Page Cleaners to start page cleaning.

- The *Page Cleaner* manages the page cleaning of the buffer pool. It accepts StartCleaning events and performs page cleaning on the dirty pages of the buffer pool. There can be more than one page cleaners in the simulator. Both the Buffer Pool Manager and the Page Cleaner can access *Buffer Pool Pages* which is the data structure holding all the pages of the buffer pool.
- An *Agent* represents a client that sends out requests to the Buffer Pool. Requests belong to each client are organized into a separated *Trace File*. Therefore, the Agent simply reads a record (either fix or unfix) from its trace file and returns it to the Buffer Pool Manager. The number of Agents is the same as the number of clients when running TPC-C.
- The *Disk* accepts disk I/O events (Read, SyncWrite, and AsyncWrite) and return IOFinish events. Because the track and sector information of disk pages is not available from the trace, it is assumed that a disk has constant read and write time. Unprocessed requests to a disk are queued and served in a first-in-first-out order. When there are multiple disks in the simulator, different disks can perform reads and writes simultaneously.

A timestamp is associated with each event. All events will be sent to the event queue and

sorted by their timestamps. The *Event Dispatcher* selects the event with the minimum timestamp and sends it to the corresponding component. This is a typical event flow when fixing a page: a fix event is read from the Agent and sent to the Buffer Pool Manager. If the page is in the buffer pool, the fix finishes, and the Buffer Pool Manager sends a NextRecord event to the Agent for the next record. If the page is not in the buffer pool, and a clean page is found for the replacement, a Read event will be sent to the Disk, and when the Buffer Pool Manager receives an IOFinish event, it will send a NextRecord event to the Agent asking for the next record.

All times in the simulator are relative, with the crucial activities shown in Table 3. This set of values was derived experimentally in order to give realistic behaviour. By comparing simulator results with results from our experimental DB2 test bed it was determined that roughly 40 million units of simulator time corresponds to 1 minute of real time on the DB2 test bed. This means, for example, that a disk operation on our system requires approximately 9 msec, a reasonable value for our configuration. These time parameters can be changed to reflect different configurations.

Table 3: Time Parameters in the Simulator

CPU operation time	
FixTime	20
UnfixTime	16
GatherCleanPage	14
ListReadTime	1
BufferpoolReadTime	1
ReplacementTime	4
Disk operation time	
ReadTime	6000
WriteTime	6000
AsynchronousWriteTime	6000

In real life, a buffer pool of a DBMS is usually IO bound. Therefore, the CPU subsystem is not part of the simulator in order to reduce the complexity. This could be a defect when the CPU overhead is high (e.g. the number of users is high), but it is sufficient for the current

investigations.

5.1 Performance Measures

In order to determine the behaviour of the buffer pool management algorithm, various quantities related to page activity, throughput and the I/O channel are measured. The detailed measures describe the change of the state of the buffer pool over time and activities associated with the movement of pages in and out of the buffer pool. Quantities measured include the number of dirty pages in the buffer pool, the number of pending I/O requests, the number of pending writes (synchronous writes and asynchronous writes), and the number of reads over some interval of interest. The metric *Transactions Finished per Interval* is used to measure the throughput of the buffer pool. If the time interval is one minute, the *Transactions Finished per Interval* is the same as TPM. In accordance with TPC-C practice, only the number of New-Order transactions finished is reported.

5.2 Simulator Calibration and Verification

Running the simulator with the time parameters shown in Table 3 generates the throughput curve shown in Figure 4. For comparison, results from a TPC-C benchmark run on our experimental test bed (the “real system”) are presented as well. To show the initial portion of the curve more clearly, only half the transactions (roughly 100 thousand) of the trace were used to produce this particular plot. Both the simulator and the real system complete the same number of transactions in this test. For the simulated results, throughput is expressed in terms of *Transactions Finished per Interval*. The interval used is $\frac{1}{30}$ of the run time. For the real system, the interval used is 30 seconds which is also $\frac{1}{30}$ of the run time. Because the real system and the simulator achieve a similar peak throughput as the buffer pool fills up and then fall to a similar stable value as the management algorithm kicks in, these two throughput curves are considered similar.

Other verification was performed. For example, after the system throughput becomes sta-

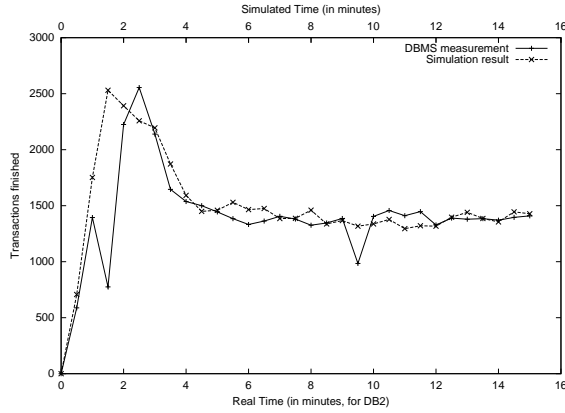


Figure 4: Simulation vs. Measurement

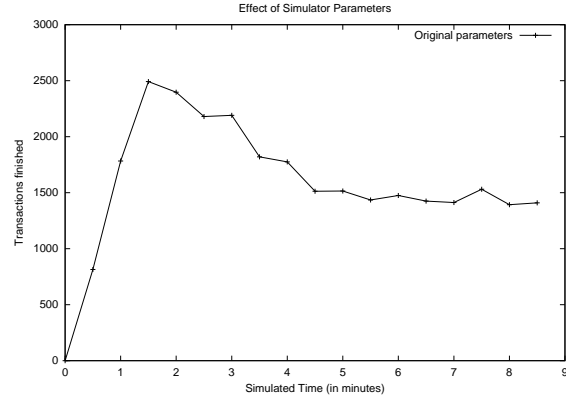


Figure 5: Original Parameters

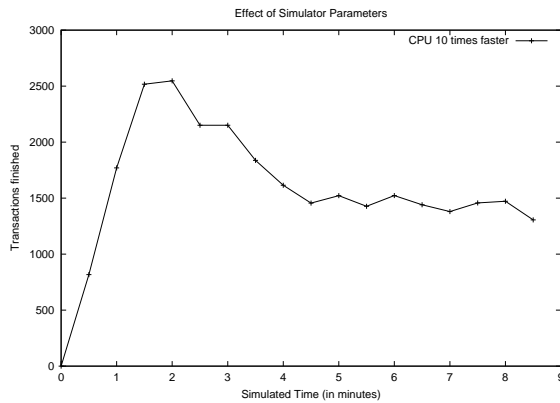


Figure 6: Faster CPU

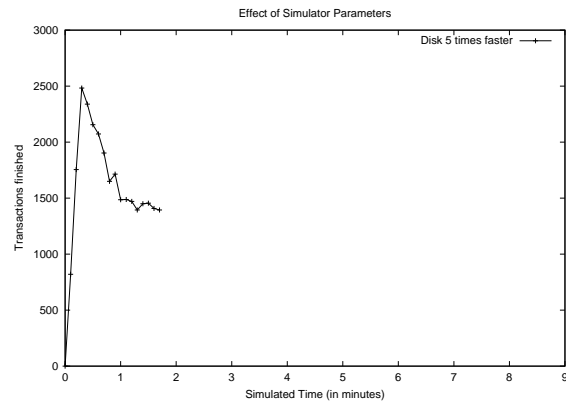


Figure 7: Faster Disk

ble, the percentage of dirty pages in the buffer pool of the simulator is 91%, which is close to the measurement result of 85%. The buffer pool hit ratio from simulation is 96.8%, which is also close to the measurement result of 96.5%.

6 Experimental Results

6.1 Effect of CPU and Disk Speed

Two experiments were performed to determine the extent to which device speeds affect performance, one with a much faster disk and the other with a much faster CPU. To focus on the change in throughput at the beginning, a shorter trace is used. Figure 5 shows the result for the original parameters. Figures 6 and 7 show the results for a faster CPU and a faster disk, respectively.

These results indicate that disk speed is more

important to DB2 than CPU speed. With a faster disk (5 times faster), the system runs faster (in this case it spends $\frac{1}{5}$ the time to complete all transactions), but with a faster CPU (10 times faster), the performance is the same as the original test. This indicates that the system is likely I/O bound.

6.2 Analysis of the Dirty Page Distribution and I/O Activities

In the experiments performed in this subsection, the default configuration of DB2 (the number of page cleaners is 2, and the dirty page threshold is 60%) was used. Figure 4 shows that the stable throughput is lower than the peak throughput under this configuration. We carried out a series of experiments to determine how various factors of the buffer pool

contribute to performance.

During the verification, we found that 87% of the pages in the buffer pool are dirty, which seems high. This motivated more experiments on the page distribution in the buffer pool. Figure 8 shows the evolution of pages in the buffer pool over the course of the simulation. Also shown on this graph is the (scaled) throughput (finished transactions per interval $\times 20$). At the beginning, all pages in the buffer pool are free pages. Both the number of dirty pages and the number of clean pages increase as time goes on. After the buffer pool is full, the number of dirty pages continues to increase, but the number of clean pages drops. At the same time, the throughput drops as well. At last, when 87% of the buffer pool pages are dirty, the system enters a steady state. The number of dirty pages at this point is much higher than it is when the buffer pool is just full, implying that there are too many dirty pages in the buffer pool. Therefore, the effect of the *dirty page threshold*, which is used to control the percentage of dirty pages in the buffer pool, was subjected to further testing.

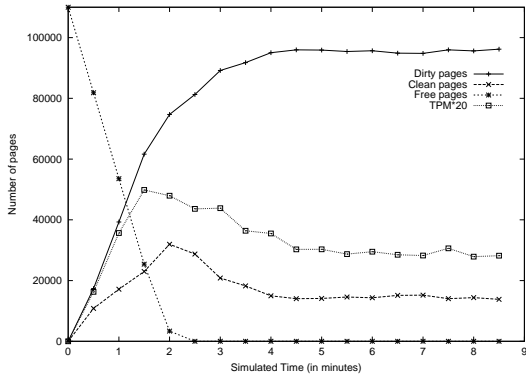


Figure 8: Pages in the Buffer Pool

Different dirty page threshold values are tested in both the simulator and DB2. The results show that the percentages of dirty pages are nearly the same regardless of the value of the threshold. This indicates that in the TPC-C workload, the threshold alone cannot effectively control the percentage of dirty pages.

To investigate the reason of this observation, the I/O activities of the buffer pool were examined in more detail. The results are shown

in Figure 9. We can see that the total I/O bandwidth is constant, and is shared by reads, synchronous writes, and asynchronous writes. The read speed is consistent with the throughput of the system. In an OLTP environment, the faster transactions are processed, the faster read requests arrive, and the faster the pages are modified.

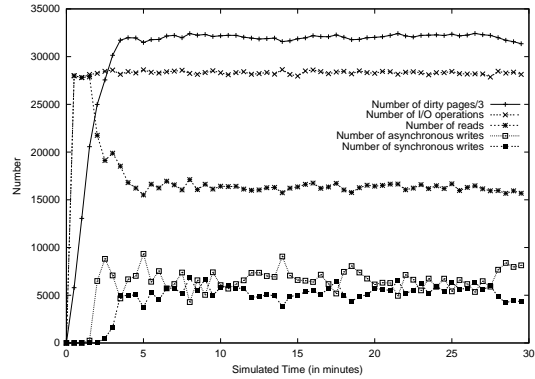


Figure 9: I/O Activities of the Buffer Pool

When the buffer pool is full, page cleaners begin to clean out dirty pages by asynchronous writes. However, asynchronous writes cannot clean out pages fast enough, and so dirty pages must be selected for replacement. Therefore, synchronous writes occur. The synchronous writes not only delay the reads directly (since a read cannot proceed before the synchronous write finishes), but also compete with other activities for I/O bandwidth. Therefore, the read speed is slowed down by the need to write in order to create space for the incoming pages. When the read speed becomes slower, the throughput drops, and dirty pages are generated more slowly. When the number of dirty pages generated by TPC-C equals the number of dirty pages cleaned by writes in the same time interval, the system enters a steady state.

As shown in Figure 9, the number of synchronous writes is high. The existence of too many synchronous writes impacts throughput. The number of asynchronous writes should be increased in order to decrease the number of synchronous writes. To do this, the aggregate page cleaning speed must be increased.

6.3 The Impact of More Page Cleaners

The aggregate page cleaning speed can be increased if more page cleaners are used. The default number of page cleaners is 2. Figure 10 shows that with 50 page cleaners, the throughput increases. Figure 11 shows the I/O activities and the number of dirty pages when the number of page cleaners is 50. There are nearly no synchronous writes left. The number of dirty pages also drops significantly.

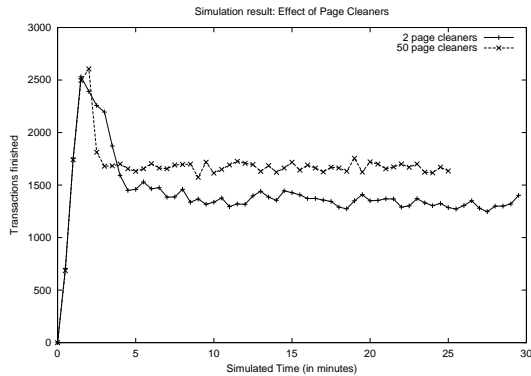


Figure 10: Effect of Number of Page Cleaners

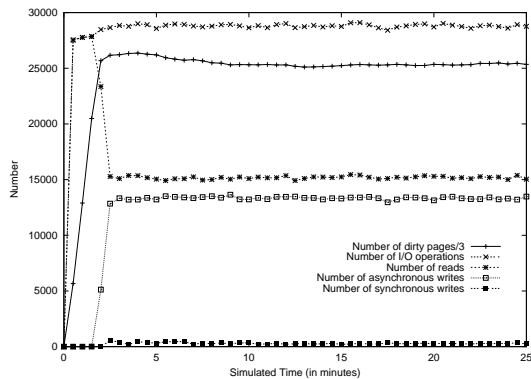


Figure 11: I/O Activities for 50 Page Cleaners

Figure 12 illustrates the throughput and dirty page percentage under different numbers of page cleaners. Increasing the number of page cleaners reduces both the percentage of dirty pages and the percentage of synchronous writes. Up to 50 page cleaners, throughput is improved with more cleaning. After that point, however, throughput drops sharply as putting

more page cleaners to work is unable to improve performance. The selection of the appropriate number of page cleaners is clearly important in tuning such a system. Similar tests on the DB2 test bed show similar results as Figure 12.

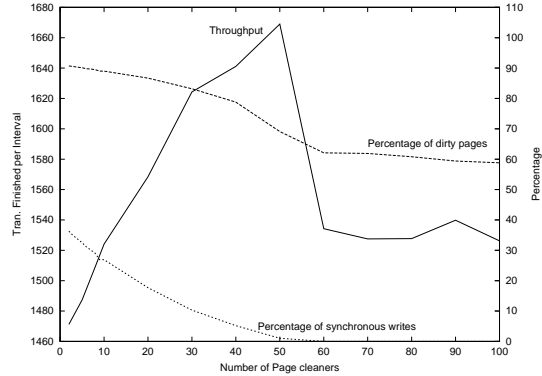


Figure 12: Effect of Multiple Page Cleaners

The page cleaning speed has an upper limit, which is achieved when all page cleaners are busy. However, the read speed has no upper limit; reads are aggressive in the TPC-C benchmark, i.e., they will consume all available disk I/O bandwidth. If the upper limit of the page cleaning speed is too low, reading can be slowed down only by synchronous writes. This is the situation with the default configuration. Increasing the number of page cleaners can increase the upper bound of the page cleaning speed. Therefore, it can decrease synchronous writes and increase throughput.

7 Self-tuning Algorithm for Page Cleaning

The above experiments show that selecting an appropriate number of page cleaners is important to the performance. However, different workloads need different number of page cleaners, and tuning such a parameter requires many experiments. To overcome this problem, a self-tuning page cleaning algorithm was developed. The objective of this algorithm is to eliminate the number of synchronous writes and limit the number of dirty pages by adjusting the disk I/O bandwidth occupied by page cleaning dynamically.

A parameter $AioP$ is introduced to control the disk I/O bandwidth occupied by page cleaning. The range of $AioP$ is $[0,1]$. Its value indicates the expected proportions of pending asynchronous writes in the total pending I/O requests. If the proportion of pending asynchronous writes is below $AioP$, the page cleaners will send out more asynchronous writes to disks. Whenever this proportion is above $AioP$, the page cleaners will stop sending any more asynchronous writes to disks, and resume cleaning only after the proportion is less than $AioP$. In this way, the $AioP$ controls the disk I/O bandwidth occupied by page cleaning. The bigger the $AioP$ value, the more disk I/O bandwidth is occupied by page cleaning. The initial value of $AioP$ can be any value between 0 and 1. Its value is adjusted based on the status of the buffer pool periodically. Before describing this algorithm, some notation is introduced:

- T : total number of pending I/O requests in the whole buffer pool.
- A : current number of pending *asynchronous* writes in the whole buffer pool.
- S : current number of pending *synchronous* writes in the whole buffer pool.
- $\delta_{d\uparrow}, \delta_{d\downarrow}, \delta_s$: scale parameters needed for this self-tuning algorithm (7.5 is used for all these parameters in the simulation experiments in this paper).
- I : the time interval that the page cleaner adjusts its parameters. The value used currently is 10 times of the disk access time.
- D_i : proportion of dirty pages in the buffer pool on the i th check interval.
- t : time.

The page cleaning control mechanism is:

1. When a page cleaner is started in response to a page cleaning event, it collects some pages for cleaning.
2. Send out n pages for cleaning, where $n = \frac{AioP * T - A}{1 - AioP}$. This lets the proportion of pending asynchronous writes reach $AioP$.
3. Whenever an asynchronous write finishes, the page cleaner performs the following check:

- (a) if $\frac{A}{T} < AioP$, send out asynchronous writes until $\frac{A}{T} = AioP$.
- (b) if $\frac{A}{T} \geq AioP$, do not send out new requests, but re-check this when the next asynchronous write finishes.

If the current interval is n , the adjustment is:

$$\begin{aligned} &\text{if } (D_n > D_{n-1}) // \text{ more dirty pages} \\ &\quad AioP \leftarrow AioP(1 + \delta_{d\uparrow} \frac{D_n - D_{n-1}}{D_{n-1}} + \delta_s S), \\ &\text{else } // \text{ less dirty pages} \\ &\quad AioP \leftarrow AioP(1 + \delta_{d\downarrow} \frac{D_n - D_{n-1}}{D_{n-1}} + \delta_s S), \\ &\text{if } AioP < 0, AioP \leftarrow 0; \\ &\text{if } AioP > 1, AioP \leftarrow 1. \end{aligned}$$

There are two terms in the adjustment. The first term, $\delta_d \frac{D_n - D_{n-1}}{D_{n-1}}$, is based on the change in dirty pages. The $\delta_{d\uparrow}$ and $\delta_{d\downarrow}$ are used to amplify this change. If the number of dirty pages drops, which means the page cleaning is too fast, $AioP$ will also drop in order to decrease the page cleaning speed; if the number of dirty pages increases, which means the page cleaning is too slow, $AioP$ will also increase in order to increase the page cleaning speed. The second term, $\delta_s S$, is used to eliminate synchronous writes. If there are many synchronous writes, $AioP$ will increase. When the number of dirty pages does not change, and there are no synchronous writes in the buffer pool, the system is in a desirable steady state, and $AioP$ does not change. If D_n and S keep changing, the $AioP$ will also keep adjusting. Because this self-tuning algorithm does not depend on any workload specific characteristics, it is expected to be useful under various workloads.

8 Test Results of the Self-Tuning Algorithm

8.1 Simulation Results

The simulation results of the self-tuning algorithm are presented in this section. Figure 13 shows the throughput of the self-tuning algorithm and the old algorithm. The throughput of the self-tuning algorithm is higher than the old algorithm. Figure 14 compares the number of dirty pages. The number of dirty

pages decreases from 91% to 69% under the self-tuning algorithm. The values of *AioP* which is scaled by a factor of 50,000 is also shown in Figure 14. The horizontal line in Figure 14 shows the maximum possible value of the *AioP*. The initial value of *AioP* is 0, and it is adjusted during the running. Figure 15 shows the I/O activities. Compare to Figure 9, the self-tuning algorithm successfully eliminates synchronous writes. The throughput comparison is shown in Table 4. The average throughput includes only the throughput after the performance drop. The self-tuning algorithm achieves a similar throughput as the hand-tuned system (the number of page cleaners is 50). Both of them have 15% higher throughput than the system under default configuration.

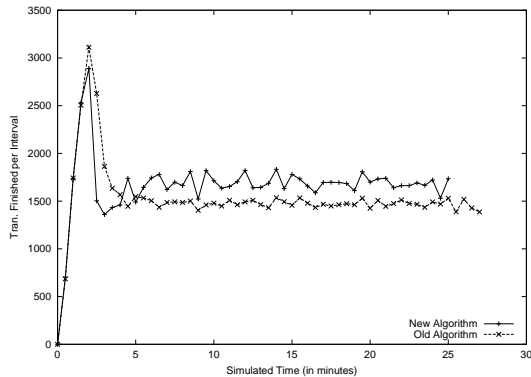


Figure 13: Throughput comparison

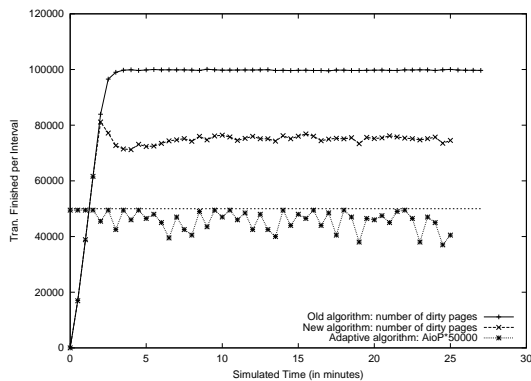


Figure 14: Dirty pages comparison and *AioP*

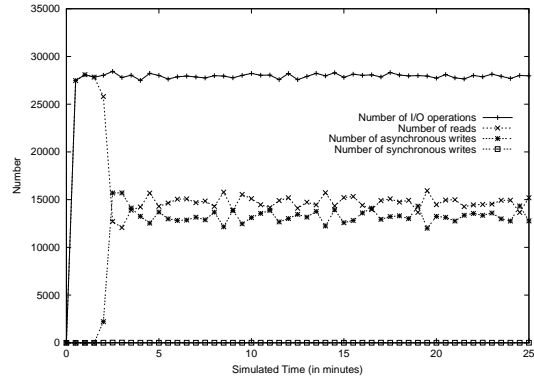


Figure 15: I/O Activities of the Self-tuning Algorithm

Table 4: Comparison of the average throughput and the C.O.V.

	Throughput	C.O.V.
Old algorithm (2 cleaners)	1471.2	0.026
Old algorithm (50 cleaners)	1669.0	0.024
Self-tuning algorithm	1691.6	0.043

8.2 Experiment Results

The self-tuning algorithm was implemented in DB2 version 7.1 and tested on systems with two different configurations. The small system contains a 50-warehouse TPC-C database spanning over 7 disks. The medium system contains a 300-warehouse TPC-C database spanning over 70 disks. The TPC-C benchmark was used as the workload. The experiment results on these two systems are presented in this section.

Figure 16 shows the throughput of the small system (50-warehouse TPC-C database) under the old algorithm with different number of page cleaners. Because the throughput with 50 page cleaners has big variance, the throughput with 16 page cleaners is selected as baseline. Table 5 shows the average TPM under different number of page cleaners and their comparison with the baseline throughput. The ratio column in the table means the ratio of the TPM to the TPM selected as baseline.

For the new algorithm, there are four param-

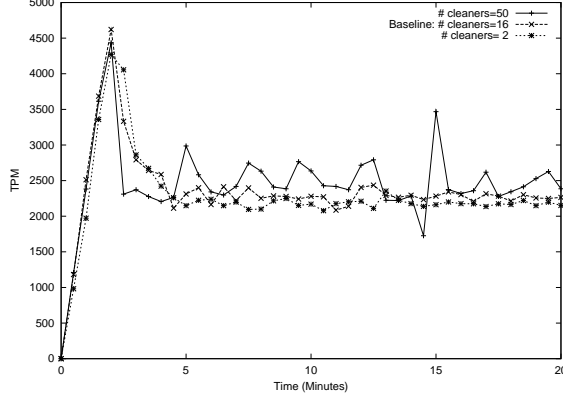


Figure 16: Small System: TPM of the Old Algorithm

Table 5: Small System: Average TPM of the Old Algorithm

# Page Cleaners	Average TPM	Ratio
50	2487.0	1.09
16 (Baseline)	2280.3	1.00
2	2177.7	0.96

eters: $\delta_{d\uparrow}$, $\delta_{d\downarrow}$, δ_s , and I . During these experiments, the I is always set to 100ms which is about 10 times of the average disk access time. It was found that there is almost no synchronous write even when $\delta_s = 0$. Therefore, δ_s is set to 0 in the experiments. The experiments showed that with more than one page cleaner, the throughput is bursty and normally not as good as the throughput with one page cleaner. Therefore, only results with one page cleaner are discussed here.

For a bigger $\delta_{d\uparrow}/\delta_{d\downarrow}$ value, the $AioP$ will be adjusted faster, and for a smaller value, the $AioP$ will be adjusted slower. However, the throughput is not sensitive to the parameters. Some results are shown in Figure 17. The parameter values are represented as a pair $(\delta_{d\uparrow}, \delta_{d\downarrow})$. The performance of the self-tuning algorithm is better than the baseline, although it is more bursty.

The throughput under the old algorithm with different number of page cleaners for the medium system (300-warehouse TPC-C database) is shown in Figure 18. Because the number of disks accessed in the system is 70, 70 page cleaners are used as one of the test case.

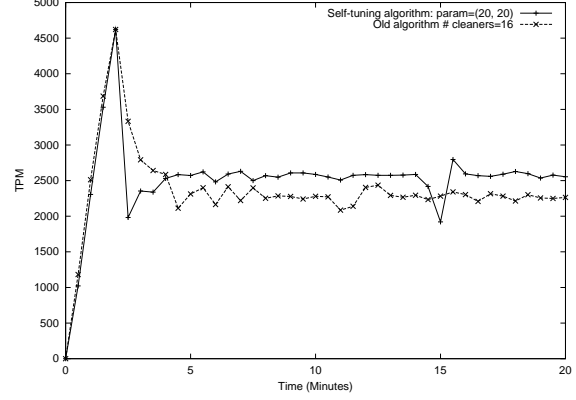


Figure 17: Small System: TPM of the New Algorithm

Table 6: Small System: Average TPM of the New Algorithm

$(\delta_{d\uparrow}, \delta_{d\downarrow})$	Average TPM	Ratio
(20, 20)	2554.5	1.12
(1, 2)	2506.3	1.10
(3, 6)	2497.8	1.10
(10, 10)	2478.3	1.09
Baseline	2280.3	1.00

The average TPM is shown in Table 7. The old algorithm with 16 page cleaners is selected as baseline.

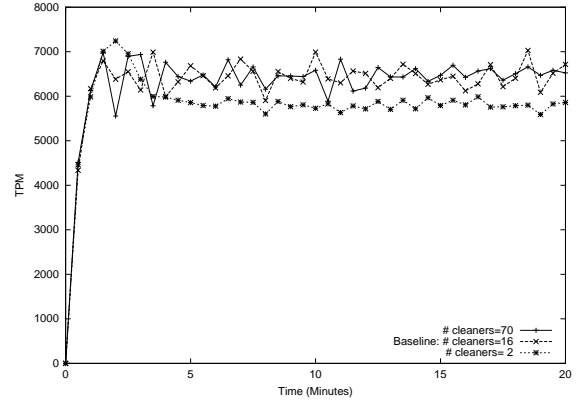


Figure 18: Medium System: TPM of the Old Algorithm

Figure 19 shows the throughput under the new algorithm. Table 8 is the average throughput. Under the self-tuning algorithm, the system throughput is comparable to the hand-

Table 7: Medium System: Average TPM

# of Page Cleaners	Average TPM	Ratio
70	6458.8	1.00
16 (Baseline)	6455.4	1.00
2	5804.1	0.90

Table 8: Medium System: Average TPM

$(\delta_{d1}, \delta_{d1})$	Average TPM	Ratio
16 (Baseline)	6455.4	1.000
(20, 20)	6383.9	0.989
(10, 10)	6373.8	0.987

tuned system.

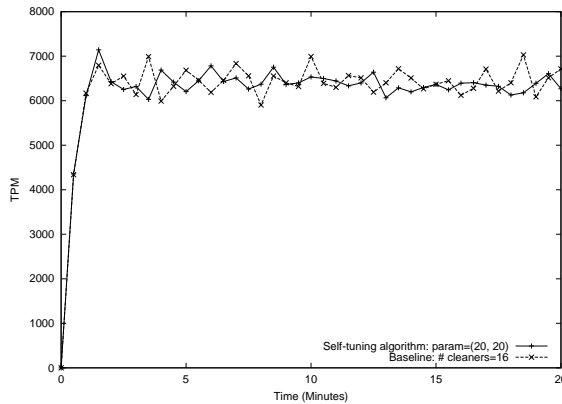


Figure 19: Medium System: TPM of the New Algorithm

9 Conclusions and Future Work

Buffer pool management is important to the performance of a DBMS. Some elements of the buffer pool management algorithm of DB2 were analyzed by trace-driven simulation, using traces captured from running the TPC-C benchmark. The effects of different parameters of the buffer pool were tested, and the I/O activities of the buffer pool were examined. Under the default configuration, although all page cleaners keep busy during the running period, dirty pages still make up more than 85% of the buffer pool. The results from further simulation experiments suggest that for this workload

the default number of page cleaners are insufficient to clean out dirty pages fast enough, and synchronous writes have to be applied as a result. Throughput is impacted negatively by the need to wait on synchronous writes. The aggregate page cleaning speed can be increased by increasing the number of page cleaners. However, the selection of an appropriate number of page cleaners is not easy. The appropriate number of page cleaners may vary when the workload changes.

A self-tuning page cleaning algorithm was developed to overcome this problem. In this self-tuning approach, only one page cleaner is needed, and the page cleaning throughput can be adjusted automatically. Simulation and measurement results showed that the system performance under this self-tuning algorithm is comparable to a hand-tuned system.

Due to the complexity of the buffer pool management algorithm, the configuration and tune of the DBMS buffer pool is a complex task. Self-tuning approaches can help the configuring and tuning of the system. The future work includes applying the self-tuning principle to other parts of the buffer pool algorithm and testing the effect in both the simulator and the real system.

References

- [1] Kurt Parick Brown. *Goal-oriented memory allocation in database management systems*. PhD thesis, University of Wisconsin-Madison, 1995.
- [2] Jen-Yao Chung, Donald Ferguson, George Wang, Christos Nikolaou, and Jim Teng. Goal oriented dynamic buffer pool management for data base systems. Technical Report TR94-0125, ICS/FORTH, Heraklion, Crete, Greece, October 1994.
- [3] Peter J. Denning. The working set model for program behavior. *Communications of the ACM*, 11(5):323–333, May 1968.
- [4] Wolfgang Effelsberg and Theo Haerder. Principles of database buffer management. *ACM Transactions on Database Systems*, 9(4):560–595, December 1984.

- [5] Wolfgang Effelsberg and Mary E. S. Loomis. Logical, internal, and physical reference behavior in CODASYL database systems. *ACM Transactions on Database Systems*, 9(2):187–213, June 1984.
- [6] Christos Faloutsos, Raymond T. Ng, and Timos K. Sellis. Flexible and adaptable buffer management techniques for database management systems. *IEEE Transactions on Computers*, 44(4):546–560, April 1995.
- [7] John P. Kearns and Samuel DeFazio. Locality of reference in hierarchical database systems. *IEEE Transactions on Software Engineering (SE)*, Vol. SE-9, (2):128–134, March 1983.
- [8] John P. Kearns and Samuel Defazio. Diversity in database reference behavior. In *Proc. of the ACM SIGMETRICS Int'l Conf. on Measurement and Modeling of Computer Systems*, volume 17(1) of *ACM SIGMETRICS Performance Evaluation Review*, pages 11–19, Berkeley, California, USA, May 23–26 1989. ACM Press.
- [9] TPC BenchmarkTM C. http://www.tpc.org/benchmark_specifications/TPC_C/tpc-c35.pdf, October 1999.