

Department of Computer Science
Database and Information Systems Group

Project Thesis:

**Performance Evaluation of Different Open
Source Implementations of Data Structures
and Other Algorithms in the context of a
DBMS Buffer Manager**

by **Max Fabian Gilbert***

Day of release: January 31, 2019

Abstract

Needless to say, every database management system needs to be able to manage data. The data structures used to manage those data in a database have a major influence on various characteristics (e.g. performance) of a database management system and therefore, the usage of specific data structures (e.g. B-tree indexes) and even some implementation details of those are very important decision in DBMS design.

But for correct and performant operation, a DBMS needs to manage various kinds of meta data as well. Some of those meta data needs to be persistent (e.g. the catalog of a relational DBMS) but some can also be non-persistent. Because of the non-persistence of data managed by the buffer management of a DB, the meta data required for the buffer manager are also usually non-persistent. The data structures used to manage those meta data are—unlike the data structures used to manage the data—more an implementation than a design decision. For some kinds of those meta data, it's—due to the non-criticality of the specific meta data management—even reasonable to use data structures provided by the used programming language even though there might be more performant data structures for the purpose. But more performant implementations for most of those data structures don't need to be implemented specifically for one project, there are many different implementations available in open source and proprietary libraries.

This work is a performance evaluation of various MPMC

This page intentionally left blank.

Contents

1	Buffer Frame Free List	1
1.1	Purpose	1
1.2	Compared Queue Implementations	1
1.2.1	Boost Lock-Free Queue with variable size	2
1.2.2	Boost Lock-Free Queue with fixed size	2
1.2.3	CDS Basket Lock-Free Queue	2
1.2.4	CDS Flat-Combining Lock-Free Queue	3
1.2.5	CDS Michael & Scott Lock-Free Queue	4
1.2.6	CDS Variation of Michael & Scott Lock-Free Queue	4
1.2.7	CDS Michael & Scott Blocking Queue with Fine-Grained Locking	4
1.2.8	CDS Ladan-Mozes & Shavit Optimistic Queue	5
1.2.9	CDS Segmented Queue	5
1.2.10	CDS Vyukov's MPMC Bounded Queue	6
1.2.11	Folly MPMC Queue	6
1.2.12	Dmitry Vyukov's Bounded MPMC Queue	6
1.2.13	Gavin Lambert's MPMC Bounded Lock-Free Queue	7
1.2.14	moodycamel::ConcurrentQueue	7
1.2.15	Matt Stump's Bounded MPMC Queue	7
1.2.16	Erik Rigtorp's Bounded MPMC Queue	8
1.2.17	Threading Building Blocks Concurrent Queue	8
1.2.18	Threading Building Blocks Bounded Concurrent Dual Queue	8
1.3	Performance Evaluation	8
1.3.1	Micro Benchmark	8
1.3.2	Used Versions of the Libraries and Queue Implementations	9
1.3.3	Configuration of the Used System	10
1.4	Conclusion	10

Bibliography

11

1 Buffer Frame Free List

1.1 Purpose

A buffer manager is required for every disk-based DBMS. A disk-based DBMS stores the pages of a database on secondary storage but to read and write pages, they are required to be in memory.

This feature is provided by the buffer pool management by managing the currently used subset of the database pages in buffer frames located in memory. A buffer frame is a portion of memory that can hold one database page and each of those frames got a frame index as identifier.

During operation, database pages are dynamically fetched from the database into buffer frames. Once a page is not required anymore, it might be evicted from the buffer pool freeing a buffer frame.

Due to the fact that pages are only allowed to be fetched into free buffer frames, the buffer manager needs to know all the free buffer frames. Therefore, a free list for the buffer frames is required.

1.2 Compared Queue Implementations

To ease implementation of page eviction strategies like CLOCK, a free list should use a FIFO data structure like a queue. Therefore the buffer frame freed first is (re-)used first as well.

Almost every state-of-the-art DBMS support multithreading and therefore, there are usually multiple threads concurrently fetching pages into the buffer pool and evicting pages from the buffer pool. Following this, a buffer frame free list has to support thread-safe functions to push frame indexes to the free list and to pop frame indexes from it. Queues providing those thread-safe access functions are usually called multi-producer (add frame indexes) multi-consumer (retrieve/remove frame indexes) queues

(MPMC queues).

An approximate number of buffer indexes in the free list must also be provided by any free list implementation to support the eviction of pages once there are only a few free buffer frames left. Thread-safe access to this number is desirable but not absolutely required.

1.2.1 Boost Lock-Free Queue with variable size

The famous *Boost C++ Libraries*¹ offer a lock-free unbounded MPMC queue² in the library `Boost.Lockfree`³. Like many other non-blocking thread-safe data structures, this MPMC queue uses atomic operations instead of locks or mutexes. To support queues of dynamically changing sizes, this queue implementation also uses a free list for the dynamic memory management internally.

This data structure does not offer the number of contained elements and therefore, an approximate number of buffer indexes in the free list needs to be managed outside.

1.2.2 Boost Lock-Free Queue with fixed size

This data structure is identical to the data structure in Subsection 1.2.1 but does not use dynamic memory management internally—it is a bounded queue. Therefore, the capacity of the queue (i.e. the maximum number of buffer frames of the buffer pool) needs to be specified beforehand which allows the usage of a fixed-size array instead of dynamically allocated nodes.

1.2.3 CDS Basket Lock-Free Queue

Besides other concurrent data structures, the *Concurrent Data Structures C++ library*⁴ offers many different thread-safe queue implementations. The

¹<https://www.boost.org/>

²<https://www.boost.org/doc/libs/release/doc/html/boost/lockfree/queue.html>

³<https://www.boost.org/doc/libs/release/doc/html/lockfree.html>

⁴<https://github.com/khizmax/libcds>

1.2 Compared Queue Implementations

unbounded *Basket Lock-Free Queue*⁵ is based on the algorithm proposed by M. Hoffman, O. Shalev and N. Shavit in [HSS07].

Internally, this queue does not use an absolute FIFO order. Instead, it puts concurrently enqueued elements into one “basket” of elements. The elements within one basket are not specifically ordered but the different “baskets” used over time are ordered according to FIFO. Therefore, the dequeue operation just dequeues one of the elements in the oldest “basket”. The dynamic memory management uses a garbage collector to deallocate emptied “baskets”.

1.2.4 CDS Flat-Combining Lock-Free Queue

The *Concurrent Data Structures C++* library does also offer an unbounded thread-safe queue that uses Flat Combining⁶. The Flat Combining technique was proposed by D. Hendler, I. Incze, N. Shavit and M. Tzafrir in [Hen+10]. This technique is used to make any sequential data structure thread-safe—in case of the *Flat-Combining Lock-Free Queue*, the `std::queue`⁷ of the *C++ Standard Library*⁸ is used as base data structure.

The Flat Combining technique uses thread-local publication lists to record operations performed by those threads. A global lock is needed to be acquired to combine these thread-local publication lists into the global, sequential data structure. The thread which acquired the global lock also combines the publication lists of all other threads reducing the locking overhead. The returned value of each operation executed during the combining is stored into the respective publication list together with the global combining pass number. A thread with a non-empty publication list that cannot acquire the global lock needs to wait till the combining thread updated its publication list.

⁵http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_basket_queue.html

⁶http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_f_c_queue.html

⁷<https://en.cppreference.com/w/cpp/container/queue>

⁸<https://en.cppreference.com/w/cpp>

1.2.5 CDS Michael & Scott Lock-Free Queue

Another unbounded lock-free queue implementation offered by the *Concurrent Data Structures* C++ library is based on the famous Michael & Scott lock-free queue algorithm⁹ which was proposed by M. Michael and M. Scott in [MS96].

The Michael & Scott lock-free queue basically uses compare-and-swap (CAS) operations on the tail of the queue to synchronize enqueue operations. If a thread reads a NULL value as next element after the queue's tail, it swaps this value atomically with the value enqueued by this thread. Afterwards it adjusts the tail pointer. If a thread does not read the NULL value there during the CAS operation, another thread has not already adjusted the tail pointer and this thread needs to retry its enqueue operation with the new tail pointer. The dequeue operation is implemented similarly. The memory occupied by already dequeued elements is deallocated using a garbage collector provided by the library.

1.2.6 CDS Variation of Michael & Scott Lock-Free Queue

The *Concurrent Data Structures* C++ library also offers an optimized variation of the Michael & Scott unbounded lock-free queue algorithm¹⁰ which is based on the works of S. Doherty, L. Groves, V. Luchangco and M. Moir in [Doh+04].

This optimization of the Michael & Scott lock-free queue optimizes the dequeue operation to only read the tail pointer once.

1.2.7 CDS Michael & Scott Blocking Queue with Fine-Grained Locking

M. Michael and M. Scott did also propose a blocking queue algorithm in [MS96]. This unbounded blocking queue implementation¹¹ is also offered

⁹http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_m_s_queue.html

¹⁰http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_moir_queue.html

¹¹http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_r_w_queue.html

1.2 Compared Queue Implementations

by the *Concurrent Data Structures C++* library.

This blocking queue algorithm uses one read and one write lock protecting the head and tail of the queue. Therefore, only one thread at a time can enqueue and only one thread at a time can dequeue elements. The deallocation of memory during dequeuing is done by the dequeuing thread instead of relying on a garbage collector.

1.2.8 CDS Ladan-Mozes & Shavit Optimistic Queue

The *Concurrent Data Structures C++* library also offers an unbounded optimistic queue implementation¹² which is based on an algorithm proposed by E. Ladan-Mozes and N. Shavit in [LS04].

Instead of using expensive CAS operations on a singly-linked list (like in the Michael & Scott lock-free queue), this algorithm uses a doubly-linked list with the possibility to detect and fix inconsistent enqueue and dequeue operations. Deallocation of memory is done using a garbage collector.

1.2.9 CDS Segmented Queue

The unbounded segmented queue implementation¹³ of the *Concurrent Data Structures C++* library is based on an algorithm proposed by Y. Afek, G. Korland and E. Yanovsky in [AKY10].

This thread-safe queue algorithm is very similar to the basket lock-free queue. It also uses a relaxed FIFO order by ordering segments containing multiple elements instead of single elements. A thread enqueueing or dequeuing elements into the tail segment or from the head segment selects one of the slots inside the segment randomly. CAS operations are used to atomically enqueue or dequeue an element from a slot. If the CAS fails, another slot is taken randomly. The size of each segment—which can be selected (8 was used for the performance evaluation in Section 1.3)—determines the relaxiness of the FIFO order. Deallocation of emptied segments is done by a garbage collector.

¹²http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_optimistic_queue.html

¹³http://libcdfs.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_segmented_queue.html

1.2.10 CDS Vyukov's MPMC Bounded Queue

The last thread-safe queue implementation¹⁴ provided by the *Concurrent Data Structures* C++ library is bounded and was developed by D. Vyukov¹⁵. The queue implementation in Subsection 1.2.12 is his original implementation.

Vyukov's thread-safe queue implementation is very similar to Michael & Scott blocking queue with fine-grained locking from Subsection 1.2.7 but instead of using mutexes as locks, his implementation uses atomic read-modify-write (**RMW**) operations. This results in a cost of basically one CAS operation per enqueue/dequeue operation.

1.2.11 Folly MPMC Queue

Facebook's open source library *Folly*¹⁶ provides a bounded lock-free queue implementation. An unbounded queue is also provided but due to the typically lower performance of unbounded queues, it is not evaluated in Section 1.3.

Folly's MPMC queue uses a ticket dispenser system to give a thread access to one of the single-element queues used. Those ticket dispensers for the head and tail of the queue use atomic increment operations which are supposed to be more robust to contention than CAS operations used e.g. in the Michael & Scott lock-free queue.

1.2.12 Dmitry Vyukov's Bounded MPMC Queue

This¹⁷ is Vyukov's original implementation of his bounded thread-safe MPMC queue.

¹⁴http://libcdis.sourceforge.net/doc/cds-api/classcds_1_1container_1_1_vyukov_m_p_m_c_cycle_queue.html

¹⁵<http://www.1024cores.net/home/lock-free-algorithms/queues/bounded-mpmc-queue>

¹⁶<https://github.com/facebook/folly>

¹⁷<http://www.1024cores.net/home/lock-free-algorithms/queues/bounded-mpmc-queue>

1.2 Compared Queue Implementations

1.2.13 Gavin Lambert's MPMC Bounded Lock-Free Queue

This¹⁸ is another implementation of Vyukov's thread-safe queue made by [Gavin Lambert].

1.2.14 `moodycamel::ConcurrentQueue`

This lock-free queue implementation¹⁹ is either unbounded or bounded depending on the used enqueueing functions and on the optional preallocation of memory (bounded behavior is used during the performance evaluation in Section 1.3). A blocking queue implementation provided by the same library is not evaluated in Section 1.3 because it is just a wrapper around the non-blocking version adding additional overhead in low-contention workloads (like the free list).

Internally, this queue implementation uses one SPMC (single producer/multiple consumer) queue per thread. Each thread enqueues elements only into its thread-local SPMC queue. When a thread tries to dequeue an element, it checks SPMC queues for emptiness until it finds one containing elements. It then dequeues one element from the SPMC queue. Therefore, this thread-safe queue does not maintain the order of elements enqueued by different threads.

Due to the implementation using multiple SPMC queues, this queue implementation should only be used as a buffer frame free list when there is exactly one thread evicting pages from the buffer pool—and therefore, enqueueing buffer frame indexes of emptied buffer frames.

1.2.15 Matt Stump's Bounded MPMC Queue

This²⁰ is another implementation of Vyukov's thread-safe queue made by Matt Stump.

¹⁸<https://gist.github.com/uecasmb547db812ae4bba39bb1bd0443801507>

¹⁹<https://github.com/cameron314/concurrentqueue>

²⁰<https://github.com/mstump/queues>

1.2.16 Erik Rigtorp's Bounded MPMC Queue

The bounded lock-free queue²¹ of Erik Rigtorp uses a ticket dispenser system similar to the one of *Folly*'s MPMC queue from Subsection 1.2.11.

1.2.17 Threading Building Blocks Concurrent Queue

The *Threading Building Blocks* library²² is an open source library originally developed by Intel®. The first thread-safe queue implementation²³ of this library is unbounded and non-blocking.

Internally, this queue implementation uses multiple lock-based micro queues to allow concurrent enqueue/dequeue executions. Therefore, the guarantees of this queue is similar to those of the `moodycamel::ConcurrentQueue` from Subsection 1.2.14.

1.2.18 Threading Building Blocks Bounded Concurrent Dual Queue

The other thread-safe queue implementation²⁴ of the *Threading Building Blocks* library is unbounded and partially non-blocking.

This queue implementation is almost identical to the other one of the Threading Building Blocks library but it does allow the limitation of the capacity. An enqueueing operation has to wait if the queue is already full according to the specified capacity.

1.3 Performance Evaluation

1.3.1 Micro Benchmark

The used micro benchmark simulates a high contented free list. The number of working threads, the number of iterations (either the fetching of a page into a free buffer frame or the eviction of a batch of pages) per thread and the batch size of buffer frames to be freed at once can be varied. It

²¹<https://github.com/rigtorp/MPMCQueue>

²²<https://www.threadingbuildingblocks.org/>

²³<https://software.intel.com/en-us/node/506200>

²⁴<https://software.intel.com/en-us/node/506201>

1.3 Performance Evaluation

does not simulate a complete buffer pool—there is only the free list with operations to enqueue and dequeue buffer frame indexes. Each working thread performs the following operations per iterations:

- If the free list is not empty:
 - Retrieve a buffer frame index from the free list.
 - Mark the retrieved buffer frame used.
- If the free list is empty:
 - While the free list is smaller than the batch eviction size:
 - * Select a random buffer frame index using a fast random numbers generator.
 - * If this buffer frame index is marked used:
 - Mark the selected buffer frame index unused.
 - Add the selected buffer frame index to the free list.

1.3.2 Used Versions of the Libraries and Queue Implementations

- *Boost C++ Libraries* 1.58
- *Concurrent Data Structures* C++ library 2.3.3²⁵
- *Folly* a15fcb1e76²⁶
- Dmitry Vyukov’s Original MPMC Queue as of September 2017
- Gavin Lambert’s MPMC Queue as of September 2017²⁷
- moodycamel::ConcurrentQueue 9f9c4e0cf4²⁸
- Matt Stump’s MPMC Queue 319c253d68²⁹

²⁵<https://github.com/khizmax/libcds/tree/5fc87a172bd82f8a7040b8b83f32ce0e635e82ea>

²⁶<https://github.com/facebook/folly/tree/a15fcb1e76444f7d464b263ad37bf3b5fbfdf33e>

²⁷<https://gist.github.com/uecasm/b547db812ae4bba39bb1bd0443801507/e40906811cb14118d328c353250559fe359f3ba7>

²⁸<https://github.com/cameron314/concurrentqueue/tree/9f9c4e0cf400bcc5c27a041e524f04e950736b25>

²⁹<https://github.com/mstump/queues/tree/319c253d68f14ac9593c3727d1597a87af73c99b>

- Erik Rigtorp's MPMC Queue 57366e41f3³⁰
- Intel® Threading Building Blocks 2017 Update 7

1.3.3 Configuration of the Used System

- **CPU:** 2× *Intel® Xeon® Processor E5420 @2.50GHz* released late 2007
- **Main Memory:** 8 × 4GB = 32GB of DDR2-SDRAM @667MHz
- **Storage:** 3 × 256GB *Samsung SSD 840 PRO Series* released mid 2012
The following data are stored on separate SSD:
 - database file of *Zero* (--sm_dbfile)
 - log directory of *Zero* (--sm_logdir)
 - log file of the buffer log for *Zero* (--sm_fix_stats_file)
 - *XtraDB* data file of *MariaDB* (datadir)
- **OS:** *Ubuntu 15.04*
- **Kernel:** *Linux 3.19.0-15-generic*
- **C++-Compiler:** *GCC (GNU Compiler Collection) 5.4.1*

1.4 Conclusion

³⁰<https://github.com/rigtorp/MPMCQueue/tree/57366e41f3f48316f175c2e704795f519a92e1d5>

Bibliography

- [AKY10] Yehuda Afek, Guy Korland, and Eitan Yanovsky. “Quasi-Linearizability: Relaxed Consistency for Improved Concurrency”. In: *Principles of Distributed Systems*. Lecture Notes in Computer Science (6490 2010): *14th International Conference, OPODIS 2010 Tozeur, Tunisia, December 14-17, 2010 Proceedings*. Ed. by Chenyang Lu, Toshimitsu Masuzawa, and Mohamed Mosbah. Found. by Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen, pp. 395–410. DOI: 10.1007/978-3-642-17653-1_29. URL: <http://mcg.cs.tau.ac.il/papers/opodis2010-quasi.pdf> (visited on Jan. 24, 2019).
- [Doh+04] Simon Doherty et al. “Formal Verification of a Practical Lock-Free Queue Algorithm”. In: *Formal Techniques for Networked and Distributed Systems – FORTE 2004*. Lecture Notes in Computer Science (3235 2004): *24th IFIP WG 6.1 International Conference, Madrid Spain, September 27-30, 2004 Proceedings*. Ed. by David de Frutos-Escrig and Manuel Núñez. Found. by Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen, pp. 97–114. DOI: 10.1007/978-3-540-30232-2_7.
- [Hen+10] Danny Hendler et al. “Flat Combining and the Synchronization-Parallelism Tradeoff”. In: (2010): *SPAA ’10: Proceedings of the Twenty-second Annual ACM Symposium on Parallelism in Algorithms and Architectures*. Ed. by Friedhelm Meyer auf der Heide and Cynthia Phillips, pp. 355–364. DOI: 10.1145/1810479.1810540. URL: <https://www.cs.bgu.ac.il/~hendlerd/papers/flat-combining.pdf> (visited on Jan. 24, 2019).
- [HSS07] Moshe Hoffman, Ori Shalev, and Nir Shavit. “The Baskets Queue”. In: *Principles of Distributed Systems*. Lecture Notes in Computer Science (4878 2007): *11th International Conference*,

- OPODIS 2007 Guadeloupe, French West Indies, December 17-20, 2007 Proceedings*. Ed. by Eduardo Tovar, Philippas Tsigas, and Hacène Fouchal. Found. by Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen, pp. 401–414. DOI: 10.1007/978-3-540-77096-1_29. URL: <https://people.csail.mit.edu/shanir/publications/Baskets%20Queue.pdf> (visited on Jan. 24, 2019).
- [LS04] Edya Ladan-Mozes and Nir Shavit. “An Optimistic Approach to Lock-Free FIFO Queues”. In: *Distributed Computing. Lecture Notes in Computer Science (3274 2004): 18th International Conference, DISC 2004, Amsterdam, The Netherlands, October 4-7, 2004 Proceedings*. Ed. by Rachid Guerraoui. Found. by Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen, pp. 117–131. DOI: 10.1007/978-3-540-30186-8_9. URL: http://people.csail.mit.edu/shanir/publications/FIFO_Queues.pdf (visited on Jan. 24, 2019).
- [MS96] Maged M. Michael and Michael L. Scott. “Simple, Fast, and Practical Non-Blocking and Blocking Concurrent Queue Algorithms”. In: (1996): *PODC ’96: Proceedings of the Fifteenth Annual ACM Symposium on Principles of Distributed Computing*. Ed. by James E. Burns and Yoram Moses, pp. 267–275. DOI: 10.1145/248052.24810629. URL: http://www.cs.rochester.edu/~scott/papers/1996_PODC_queues.pdf (visited on Jan. 24, 2019).