

**VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING**



SPECIALIZED PROJECT

**STUDYING AND DEVELOPING
NONBLOCKING DISTRIBUTED MPSC QUEUES**

Major: Computer Science

THESIS COMMITTEE: 6 - ĐỒ ÁN CHUYÊN NGÀNH

SUPERVISORS: DR. DIỆP THANH ĐĂNG

ASSOC. PROF. DR. THOẠI NAM

REVIEWER:

—000—

STUDENT: ĐỖ NGUYỄN AN HUY - 2110193

HCMC, 08/2025

Disclaimers

I affirm that this specialized project is the product of my original research and experimentation. Any references, resources, results which this project is based on or a derivative work of have been given due citations and properly listed in the references section. All original contents presented are the culmination of my dedication and perserverance under the close guidance of my supervisors, Mr. Thoại Nam and Mr. Diệp Thanh Đăng, from the Faculty of Computer Science and Engineering, Ho Chi Minh City University of Technology. I take full responsibility for the accuracy and authenticity of this document. Any misinformation, copyright infringement or plagiarism shall be faced with serious punishment.

Acknowledgements

This thesis is the culmination of joint efforts coming from not only myself, but also my professors, my family, my friends and other teachers of Ho Chi Minh University of Technology.

I want to first acknowledge my university, Ho Chi Minh University of Technology. Throughout my four years of pursuing education here, I have built a strong theoretical foundation and gained various practical experiences. These all lend themselves well to the completion of this thesis. However, I have to say that the one person who plays the cornerstone in my academic foundation must be Mr. Diệp Thanh Đăng. He has been the constant supervisor and reliable consultant through this capstone project, right from its inception, throughout its nurturing until the very end. The project couldn't have reached this stage of maturity without Mr. Đăng. Therefore, he has the right to share whatever yields this project has to offer. I know it is getting informal, but if you ever read this, I want you to know that you surely have become one of the most important people in my life. You are the right blend of intellect, wisdom, organization, empathy, open-mindedness that I always seek in an academic companion, which surely secures you a special place in my heart. Whatever the outcome of what's about to happen, whoever I become in the future, I won't ever forget you and our magical moments when we were brewing up this project. I want to credit you for anything I achieve in front of just anyone, including myself. I have always craved a companionship like Friedrich Engels and Karl Marx; maybe I have come quite close to it.

I want to thank Mr. Thoại Nam for providing all the facilities and infrastructure that this project necessitates. He essentially funded all of our idea-brewing grounds, where magic happens.

I also want to give my family the sincerest thanks for their emotional and financial support, without which I couldn't have wholeheartedly followed my research till the end.

Last but not least important, I want to thank my closest friends for their informal but ever-constant check-ups to make sure I didn't miss the timeline for this specialized project, which I usually do not have the mental capacity for.

Contents

Chapter I Introduction	9
1.1 Motivation	9
1.2 Objective	10
1.3 Scope	11
1.4 Research question	11
1.5 Thesis overview	11
1.6 Structure	12
Chapter II Background	14
2.1 Irregular applications	14
2.1.1 Actor model as an irregular application	15
2.1.2 Fan-out/Fan-in pattern as an irregular application	16
2.2 MPSC queue	17
2.3 Correctness condition of concurrent algorithms	17
2.4 Progress guarantee of concurrent algorithms	19
2.4.1 Blocking algorithms	19
2.4.2 Non-blocking algorithms	20
2.4.2.1 Lock-free algorithms	20
2.4.2.2 Wait-free algorithms	21
2.5 Popular atomic instructions in designing non-blocking algorithms	21
2.5.1 Fetch-and-add (FAA)	21
2.5.2 Compare-and-swap (CAS)	22
2.5.3 Load-link/Store-conditional (LL/SC)	23
2.6 Common issues when designing non-blocking algorithms	23
2.6.1 ABA problem	24
2.6.2 Safe memory reclamation problem	25
2.7 MPI-3 - A popular distributed programming library interface specification ..	26
2.7.1 MPI-3 RMA	26
2.7.2 MPI-RMA communication operations	26
2.7.3 MPI-RMA synchronization	27
2.8 Pure MPI - A porting approach of shared memory algorithms to distributed algorithms	28
2.9 BCL CoreX	30
Chapter III Related works	32
3.1 Non-blocking shared-memory MPSC queues	32
3.1.1 LTQueue	32
3.1.2 DQueue	34
3.1.3 WRLQueue	36
3.1.4 Jiffy	37
3.1.5 Remarks	38
3.2 Distributed MPSC queues	38
Chapter IV Distributed MPSC queues	41

4.1 A simple baseline distributed SPSC	42
4.2 dLTQueue - Straightforward LTQueue adapted for distributed environment .	45
4.2.1 Overview	45
4.2.2 Data structure	46
4.2.3 Algorithm	48
4.3 Slotqueue - dLTQueue-inspired distributed MPSC queue with all constant-time operations	54
4.3.1 Overview	54
4.3.2 Data structure	55
4.3.3 Algorithm	56
Chapter V Evaluation	60
5.1 Benchmarking metrics	60
5.1.1 Throughput	60
5.1.2 Latency	60
5.2 Benchmarking baselines	61
5.3 Microbenchmark program	61
5.4 Benchmarking setup	61
5.5 Benchmarking results	62
Chapter VI Conclusion	64
Appendix A Theoretical aspects	65
A.1 Preliminaries	65
A.1.1 Linearizability	65
A.1.2 Larch interface specification of queues	66
A.1.3 Verifying linearizability	66
A.1.4 The Owicki-Gries method	66
A.1.5 ABA-safety	66
A.2 Theoretical proofs of the distributed SPSC	66
A.2.1 Correctness	66
A.2.1.1 ABA problem	67
A.2.1.2 Memory reclamation	67
A.2.1.3 Linearizability	67
A.2.2 Progress guarantee	67
A.2.3 Theoretical performance	67
A.3 Theoretical proofs of dLTQueue	68
A.3.1 Proof-specific notations	68
A.3.2 Correctness	70
A.3.2.1 ABA problem	70
A.3.2.2 Memory reclamation	70
A.3.2.3 Linearizability	70
A.3.3 Progress guarantee	70
A.3.4 Theoretical performance	70
A.4 Theoretical proofs of Slotqueue	71
A.4.1 Proof-specific notations	71

A.4.2 Correctness	73
A.4.2.1 ABA problem	73
A.4.2.2 Memory reclamation	73
A.4.2.3 Linearizability	73
A.4.3 Progress guarantee	75
A.4.4 Theoretical performance	75
References	77

List of Tables

Table 1	Specification of <code>MPI_Win_lock_all</code> and <code>MPI_Win_unlock_all</code>	29
Table 2	Summary of existing shared memory MPSC queues. The cell marked with (*) indicates that our evaluation contradicts with the authors' claims.	32
Table 3	Characteristic summary of existing distributed MPSC queues. <i>R</i> stands for remote operations and <i>L</i> stands for local operations. (*) [1] claims that it is lock-free.	39
Table 4	Characteristic summary of our proposed distributed MPSC queues. (1) <i>n</i> is the number of enqueueers. (2) <i>R</i> stands for remote operation and <i>L</i> stands for local operation . (*) The underlying SPSC is assumed to be our simple distributed SPSC in Section 4.1.	42
Table 5	Theoretical performance summary of our simple distributed SPSC. <i>R</i> means remote operations and <i>L</i> means local operations.	67
Table 6	Theoretical performance summary of <code>dLTQueue</code> . <i>R</i> means remote operations and <i>L</i> means local operations.	71
Table 7	Theoretical performance summary of <code>SlotQueue</code> . <i>R</i> means remote operations and <i>L</i> means local operations.	76

List of Images

Figure 1	An overview of this thesis.	12
Figure 2	Actor model visualization.	15
Figure 3	Fan-out/Fan-in pattern visualization.	16
Figure 4	Linearization points of method 1, method 2, method 3, method 4 happen at $t_1 < t_2 < t_3 < t_4$, therefore, their effects will be observed in this order as if we call method 1, method 2, method 3, method 4 sequentially.	18
Figure 5	Blocking algorithm: When a process is suspended, it can potentially block other processes from making further progress.	19
Figure 6	Lock-free algorithm: All the live processes together always finish in a finite amount of steps.	20
Figure 7	Wait-free algorithm: Any live process always finishes in a finite amount of steps.	21
Figure 8	ABA problem in a linked-list stack.	24
Figure 9	Unsafe memory reclamation in a LIFO stack.	25
Figure 10	An illustration of passive target communication. Dashed arrows represent synchronization (source: [2]).	28
Figure 11	An illustration of our synchronization approach in MPI RMA.	30
Figure 12	LTQueue's structure.	33
Figure 13	Intuition on how timestamp-refreshing works.	34
Figure 14	DQueue's structure.	35
Figure 15	WRLQueue's structure.	36
Figure 16	WRLQueue's dequeue operation	37
Figure 17	Jiffy's structure.	37
Figure 18	AMQueue's structure.	39
Figure 19	dLTQueue's structure.	46
Figure 20	Basic structure of Slotqueue.	55
Figure 21	Microbenchmark results for enqueue operation.	62
Figure 22	Microbenchmark results for dequeue operation.	62
Figure 23	Microbenchmark results for total throughput.	63
Figure 24	dLTQueue's structure.	69
Figure 25	Basic structure of Slotqueue.	72

Chapter I Introduction

This chapter details the motivation for our research topic: “Studying and developing non-blocking distributed MPSC queues”, based on which we set out the objectives and scope of this study. To summarize, we then come to the formulation of our research question and give a high-level overview of the thesis. We end this chapter with a brief description of the structure of the rest of this document.

1.1 Motivation

The demand for computation power has been increasing relentlessly. Increasingly complex computation problems arise and accordingly more computation power is required to solve them. Much engineering effort has been put forth toward obtaining more computation power. A popular topic in this regard is distributed computing: The combined power of clusters of commodity hardware can surpass that of a single powerful machine [3].

To harness the power of distributed systems, specialized algorithms and data structures need to be devised. Two especially important properties of distributed systems are performance and fault tolerance [4]. Therefore, the algorithms and data structures running on distributed systems need to be highly efficient and fault tolerant. Regarding efficiency, we are concerned with the algorithms’ throughput and latency, which are the two main metrics to measure performance. Considering fault tolerance, we are especially interested in the progress guarantee [5] characteristic of the algorithms. The progress guarantee criterion divides the algorithms into two groups: blocking and non-blocking. Blocking algorithms allow one faulty process to delay the other processes forever, which is not fault tolerant [6]. Non-blocking algorithms are safeguarded against this problem, exhibiting a higher degree of fault tolerance [7].

One of the algorithms that has seen applications in the distributed domain is the multi-producer, single-consumer (MPSC) queue algorithm [1]. Furthermore, there are applications and programming patterns in the shared-memory domain that can potentially see similar usage in the distributed domain, such as the actor model [8] or the fan-out fan-in pattern [9]. Although the more general multi-producer, multi-consumer (MPMC) queues suffice for the MPSC workloads, they are typically too expensive for these use cases [10], [11]. Therefore, supporting a specialized non-blocking distributed MPSC queue is still valuable.

However, currently in the literature, there is only one distributed MPSC queue, AMQueue [1]. Moreover, even though the author claims that AMQueue is non-blocking, we found that AMQueue is actually blocking (Section 3.2). This is unlike the shared-memory domain, where there are a lot more research on non-blocking MPSC queues [10], [11], [12], [13]. This apparent gap between the two domains have been bridged by some recent research to adapt non-blocking shared-memory algorithms to distributed environments [14], [15], [16], [17]. The work by [17] introduces a method for creating non-blocking distributed data structures within the partitioned global

address space (PGAS) framework, particularly targeting the Chapel programming language. However, their methodology faces a significant limitation: it relies on double-word compare-and-swap (DCAS) or 128-bit compare-and-swap (CAS) operations to prevent ABA problems, which lack support from most remote direct memory access (RDMA) hardware systems [17]. The HCL framework [16] provides a distributed data structure library built on RPC over RDMA technology. While functional, this approach demands specialized hardware capabilities from contemporary network interface cards, limiting its portability [15]. BCL Core [14] presents a highly portable solution capable of interfacing with multiple distributed programming backends including MPI, SHMEM, and GASNet-EX. However, BCL Core's architecture incorporates 128-bit pointers, creating the same RDMA hardware compatibility issues as [17]. For our research, we have selected BCL CoreX [15] and adopted its design philosophy to adapt existing shared-memory MPSC queues for distributed computing environments. BCL CoreX [15] extends the original BCL [14] framework with enhanced features that simplify the development of non-blocking distributed data structures. A key innovation in their approach is the implementation of 64-bit pointers, which are compatible with virtually all large-scale computing clusters and supported by most RDMA hardware configurations. To address ABA problems without relying on specialized instructions like DCAS, they have developed a distributed hazard pointer mechanism. This generic solution provides sufficient portability and flexibility to accommodate the adaptation of most existing non-blocking shared-memory data structures to distributed environments.

In summary, we focus on the design of efficient non-blocking distributed MPSC queues using the BCL CoreX library as the main implementation framework. The next few sections will list the objectives in more details (Section 1.2, Section 1.3), sum them up in a research question (Section 1.4) and an overview picture of the thesis (Section 1.5).

1.2 Objective

Based on what we have listed out in Section 1.1, we aim to:

- Investigate the principles underpinning the design of fault-tolerant and performant shared-memory algorithms.
- Investigate state-of-the-art shared-memory MPSC queue algorithms as case studies to support our design of distributed MPSC queue algorithms.
- Investigate existing distributed MPSC algorithms to serve as a comparison baseline.
- Model and design distributed MPSC queue algorithms using techniques from the shared-memory literature, specifically the BCL CoreX library.
- Utilize the shared-memory programming model to evaluate various theoretical aspects of distributed MPSC queue algorithms: correctness and progress guarantee.
- Model the theoretical performance of distributed MPSC queue algorithms that are designed using techniques from the shared-memory literature.

- Collect empirical results on distributed MPSC queue algorithms and discuss important factors that affect these results.

1.3 Scope

The following narrows down what we are going to investigate in the shared-memory literature and which theoretical and empirical aspects we are interested in for our distributed algorithms:

- Regarding the investigation of the design principles in the shared-memory literature, we focus on fault-tolerant and performant concurrent algorithm design using atomic operations and common problems that often arise in this area, namely, ABA problem and safe memory reclamation problem.
- Regarding the investigation of shared-memory MPSC queues currently in the literature, we focus on linearizable MPSC queues that follow strict FIFO semantics and support at least lock-free enqueue and dequeue operations.
- Regarding correctness, we concern ourselves with the linearizability correctness condition.
- Regarding fault tolerance, we concern ourselves with the concept of progress guarantee, that is, the ability of the system to continue to make forward progress despite the failure of one or more components of the system.
- Regarding algorithm prototyping, benchmarking and optimizations, we assume an MPI-3 setting.
- Regarding empirical results, we focus on performance-related metrics, e.g. throughput and latency.

1.4 Research question

Any research effort in this thesis revolves around this research question:

“How to utilize shared-memory programming principles to model and design distributed MPSC queue algorithms in a correct, fault-tolerant and performant manner?”

We further decompose this question into smaller subquestions:

1. How to model the correctness of a distributed MPSC queue algorithm?
2. Which factors contribute to the fault tolerance and performance of distributed MPSC queue algorithms?
3. Which shared-memory programming principles are relevant in modeling and designing distributed MPSC queue algorithms in a fault-tolerant and performant manner?
4. Which shared-memory programming principles need to be modified to more effectively model and design distributed MPSC queue algorithms in a fault-tolerant and performant manner?

1.5 Thesis overview

An overview of this thesis is given in Figure 1.

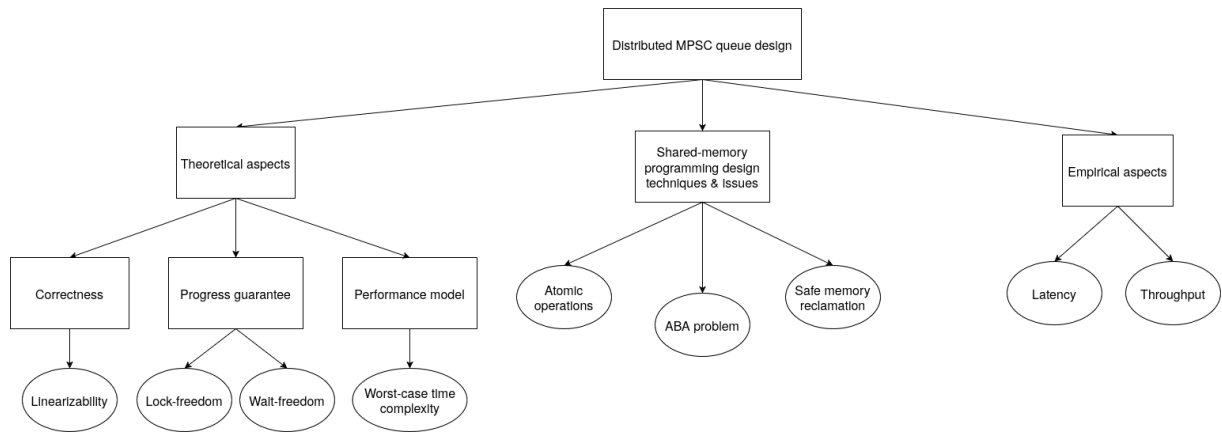


Figure 1: An overview of this thesis.

This thesis explores the shared-memory programming model to design fault-tolerant and performant concurrent algorithms using atomic operations. Traditionally, in this aspect, two notorious problems often arise: ABA problem and safe memory reclamation. We investigate the traditional techniques used in the shared-memory literature to resolve these problems and appropriately adapt them to solve similar issues when designing fault-tolerant and performant distributed MPSC queues.

This thesis contributes two new wait-free distributed MPSC queue algorithms. Theoretically, we are concerned with their correctness (linearizability), progress guarantee (lock-freedom and wait-freedom) which has an implication on their fault tolerance and their theoretical performance, which is approximated by their number of remote operations and local operations.

This thesis concludes with an empirical analysis of our novel algorithms to see if their actual behavior matches our theoretical performance model, interprets these results and discusses their implications.

1.6 Structure

The rest of this report is structured as follows:

Chapter II discusses the theoretical foundation this thesis is based on. As mentioned, this thesis investigates the principles of shared-memory programming and the existing state-of-the-art shared-memory MPSC queues. We then explore the utilities offered by MPI-3 and BCL CoreX to implement distributed algorithms modeled by shared-memory programming techniques.

Chapter III surveys the shared-memory literature for state-of-the-art queue algorithms, specifically MPSC queues. We specifically focus on non-blocking shared-memory algorithms that have the potential to be adapted efficiently for distributed environments. This chapter additionally surveys existing distributed MPSC algorithms to serve as a comparison baseline for our novel distributed MPSC queue algorithms.

Chapter IV introduces our novel distributed MPSC queue algorithms, designed using shared-memory programming techniques and inspired by the selected shared-memory

MPSC queue algorithms surveyed in Chapter III. It specifically presents our adaptation efforts of existing algorithms in the shared-memory literature to make their distributed implementations feasible and efficient.

Chapter V details our benchmarking metrics and elaborates on our benchmarking setup. We aim to demonstrate results on how well our novel MPSC queue algorithms perform, additionally compared to existing distributed MPSC queues. Finally, we discuss important factors that affect the runtime properties of distributed MPSC queue algorithms.

Chapter VI concludes what we have accomplished in this thesis and considers future possible improvements to our research.

Chapter II Background

This chapter provides various information about the terminology referenced throughout this thesis. To motivate the discussion of MPSC queues in Section 2.2, we first discuss two irregular applications in Section 2.1. Next, we decide what it means for a concurrent algorithm to be correct in Section 2.3 and the progress guarantee characteristics of concurrent algorithms in Section 2.4. From there, we decide to design linearizable non-blocking distributed MPSC queues. Therefore, we are concerned with the tools needed to design non-blocking algorithms in Section 2.5 and the issues that arise in this design process such as ABA problem and safe memory reclamation problem in Section 2.6. We finally introduce the practical libraries to help us realize non-blocking distributed MPSC queues in Section 2.7, Section 2.8 and Section 2.9.

2.1 Irregular applications

MPSC queue (Section 2.2) and its applications belong to a class called irregular applications. Designing irregular applications needs to take into account their special properties, which motivates Section 2.7, Section 2.8, Section 2.9. Therefore, before we discuss MPSC queue in Section 2.2, we explain the term “irregular application” in this section.

Irregular applications [18] are a class of programs particularly interesting in distributed computing. They are characterized by:

- Unpredictable memory access: Before the program is actually run, we cannot know which data it will need to access. We can only know that at run time.
- Data-dependent control flow: The decision of what to do next (such as which data to access next) is highly dependent on the values of the data already accessed, hence the unpredictable memory access property because we cannot statically analyze the program to know which data it will access. The control flow is inherently engraved in the data, which is not known until runtime.

Irregular applications are interesting because they demand special techniques to achieve high performance [18]. One specific challenge is that this type of application is hard to model in traditional MPI APIs using the Send/Receive interface [19]. This is specifically because using this interface requires a programmer to have already anticipated communication within pairs of processes before runtime, which is difficult with irregular applications. The introduction of MPI remote memory access (RMA) in MPI-2 and its improvement in MPI-3 has significantly improved MPI’s capability to express irregular applications comfortably [20]. This will be explained further in Section 2.7.

2.1.1 Actor model as an irregular application

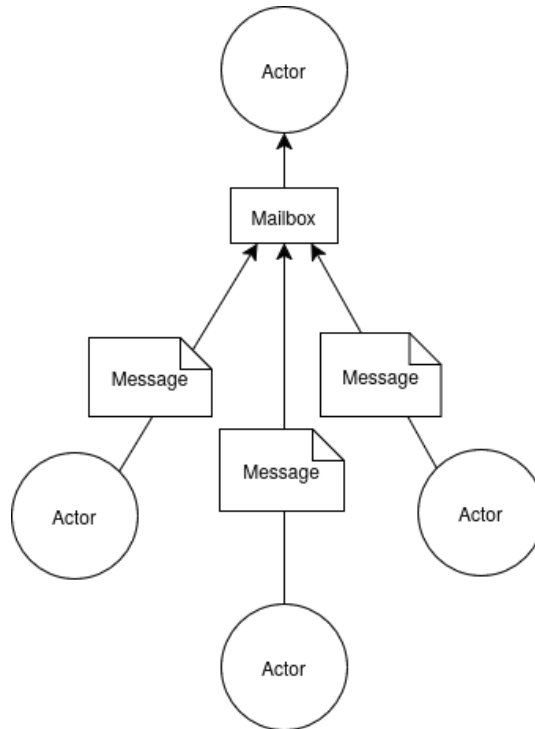


Figure 2: Actor model visualization.

The actor model [8] in actuality is a type of irregular application supported by the concurrent MPSC queue data structure.

Each actor can be a process or a compute node in the cluster, carrying out a specific responsibility in the system. From time to time, there is a need for the actors to communicate with each other. For this purpose, the actor model offers a mailbox local to each actor. This mailbox exhibits MPSC queue behavior: Other actors can send messages to the mailbox to notify the owner actor and the owner actor at their leisure repeatedly extracts messages from its mailbox. The actor model provides a simple programming model for concurrent processing.

The reasons why the actor model is an irregular application are straightforward to see:

- Unpredictable memory access: The cases in which one actor can anticipate which one of the other actors can send it a message are pretty rare and application-specific. As a general framework, in an actor model, the usual assumption is that any number of actors can try to communicate with an actor at some arbitrary time. By this nature, the communication pattern is unpredictable.
- Data-dependent control flow: If an actor A sends a message to another actor B, and when B reads this message, B decides to send another message to another actor C. As we can see, the control flow is highly engraved in the messages, or in other words, the messages drive the program flow, which can only be known at runtime.

2.1.2 Fan-out/Fan-in pattern as an irregular application

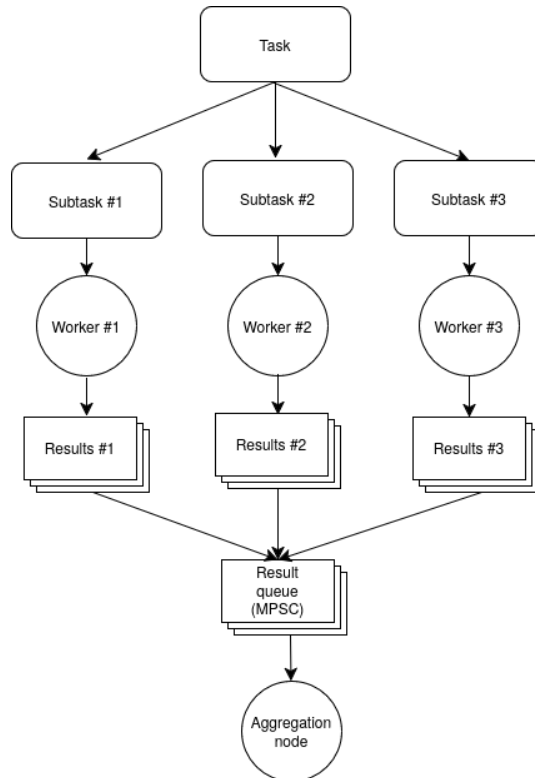


Figure 3: Fan-out/Fan-in pattern visualization.

The fan-out/fan-in pattern [9] is another type of irregular application supported by the concurrent MPSC queue data structure.

In this pattern, there is a big task that can be split into subtasks to be executed concurrently on some work nodes. In the execution process, each worker produces a result set, each enqueued back to a result queue located on an aggregation node. The aggregation node can then dequeue from this result queue to perform further processing. Clearly, this result queue exhibits MPSC behavior.

The fan-out/fan-in pattern exhibits less irregularity than the actor model, however. Usually, the worker nodes and the aggregation node are known in advance. The aggregation node can anticipate Send calls from the worker nodes. Still, there is a degree of irregularity that this pattern exhibits: How can the aggregation node know how many Send calls a worker node will issue? This is highly driven by the task and the data involved in this task, hence, we have the data-dependent control flow property. One can still statically calculate or predict how many Send calls a worker node will issue. Nevertheless, this is problem-specific. Therefore, the memory access pattern is somewhat unpredictable. Notice that if supported by a concurrent MPSC queue data structure, the fan-out/fan-in pattern is free from this burden of organizing the right amount of Send/Receive calls. Thus, combining with the MPSC queue, the fan-out/fan-in pattern becomes more general and easier to program.

We have seen the role MPSC queues play in supporting irregular applications. It is important to understand what really comprises an MPSC queue data structure.

2.2 MPSC queue

Having established the notion of irregular applications in Section 2.1, we can discuss about our design goal, distributed MPSC queue, which is an irregular application itself, in this section. The design criteria will be detailed later, in Section 2.3 and Section 2.4.

Multi-producer, single-consumer (MPSC) queue is a specialized concurrent first-in first-out (FIFO) data structure. A FIFO is a container data structure where items can be inserted into or taken out of, with the constraint that the items that are inserted earlier are taken out earlier. Hence, it is also known as the queue data structure. The process that performs item insertion into the FIFO is called the producer and the process that performs item deletion (and retrieval) is called the consumer.

In concurrent queues, multiple producers and consumers can run concurrently. One class of concurrent FIFOs is the MPSC queue, where one consumer may run in parallel with multiple producers.

The reasons we are interested in MPSC queues instead of the more general multi-producer, multi-consumer (MPMC) queue data structures are that (1) high-performance and high-scalability MPSC queues are much simpler to design than MPMCs while (2) MPSC queues are powerful enough to solve certain problems, as demonstrated in Section 2.1. The MPSC queue in actuality is an irregular application in itself:

- Unpredictable memory access: As a general data structure, the MPSC queue allows any process to enqueue and dequeue at any time. By nature, its memory access pattern is unpredictable.
- Data-dependent control flow: The consumer's behavior is entirely dependent on whether and which data is available in the MPSC queue. The execution paths of MPSC queues can vary, based on the queue contention i.e. some processes may back off or retry some failed operations; this scenario often arises in lock-free data structures.

As an implication, some irregular applications can actually “push” the “irregularity burden” to the distributed MPSC queue, which is already designed for high performance and fault tolerance. This provides a comfortable level of abstraction for programmers that need to deal with irregular applications.

2.3 Correctness condition of concurrent algorithms

We have established our design goal in the previous sections (Section 2.1, Section 2.2), that is MPSC queue. During this design process, we have to take into account its correctness, which is the subject of this section. The fault tolerance characteristic, although important, is less compared to correctness, and so will be deferred to Section 2.4.

Correctness of concurrent algorithms is hard to define, regarding the semantics of concurrent data structures like MPSC queues. One effort to formalize the correctness of concurrent data structures is the definition of **linearizability** [21]. A method call

on the FIFO can be visualized as an interval spanning two points in time. The starting point is called the **invocation event** and the ending point is called the **response event**. **Linearizability** informally states that each method call should appear to take effect instantaneously at some moment between its invocation event and response event [5]. The moment the method call takes effect is termed the **linearization point**. Specifically, suppose the following:

- We have n concurrent method calls m_1, m_2, \dots, m_n .
- Each method call m_i starts with the **invocation event** happening at timestamp s_i and ends with the **response event** happening at timestamp e_i . We have $s_i < e_i$ for all $1 \leq i \leq n$.
- Each method call m_i has the **linearization point** happening at timestamp l_i , so that $s_i \leq l_i \leq e_i$.

Then, linearizability means that if we have $l_1 < l_2 < \dots < l_n$, the effect of these n concurrent method calls m_1, m_2, \dots, m_n must be equivalent to calling m_1, m_2, \dots, m_n **sequentially**, one after the other in that order.

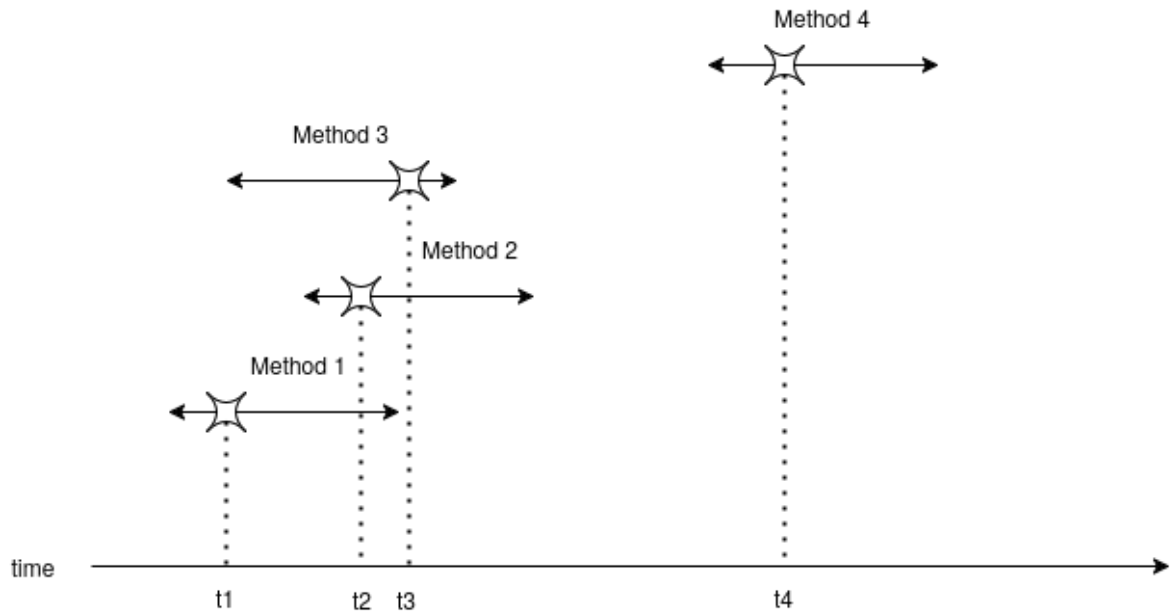


Figure 4: Linearization points of method 1, method 2, method 3, method 4 happen at $t_1 < t_2 < t_3 < t_4$, therefore, their effects will be observed in this order as if we call method 1, method 2, method 3, method 4 sequentially.

Linearizability is widely used as a correctness condition because of (1) its composability (if every component in the system is linearizable, the whole system is linearizable), which promotes modularity and ease of proof (2) its compatibility with human intuition, i.e. linearizability respects real-time order [21]. Naturally, we choose linearizability to be the only correctness condition for our algorithms.

A more formal treatment of linearizability is given in Section A.1.1.

2.4 Progress guarantee of concurrent algorithms

A correct algorithms can still be prone to faults at runtime, which varies from a process experiences an unexpected delay in its execution to a process crashes indefinitely. Therefore, fault tolerance is also an important criteria for our design goal, distributed MPSC queue (Section 2.2), besides correctness (Section 2.3). This section will introduce the concept of progress guarantee, which is highly linked with fault tolerance. The techniques to achieve fault tolerance are discussed in the next section (Section 2.5).

Progress guarantee [5] is a criterion that only arises in the context of concurrent algorithms. Informally, it is the degree of hindrance one process imposes on another process from completing its task. In the context of sequential algorithms, this is irrelevant because there is only ever one process. Progress guarantee has an implication on an algorithm's performance and fault tolerance, especially in adverse situations, as we will explain in the following sections.

2.4.1 Blocking algorithms

Many concurrent algorithms are based on locks to create mutual exclusion, in which only some processes that have acquired the locks are able to act, while the others have to wait. While lock-based algorithms are simple to read, write and verify, these algorithms are said to be **blocking**: One slow process may slow down the other faster processes, for example, if the slow process successfully acquires a lock and then the operating system (OS) decides to suspend it to schedule another one, this means until the process is awakened, the other processes that contend for the lock cannot continue.

Blocking is the weakest progress guarantee one algorithm can offer; it allows one process to impose arbitrary impedance to any other processes, as shown in Figure 5.

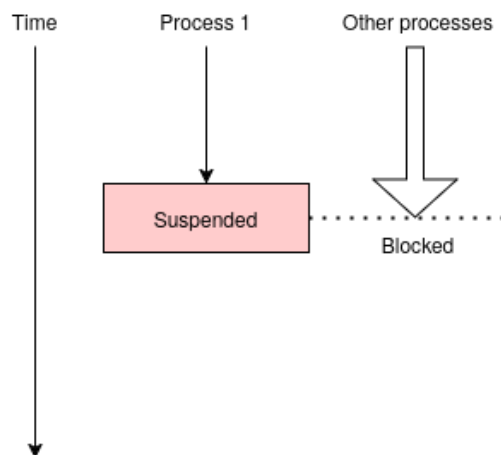


Figure 5: Blocking algorithm: When a process is suspended, it can potentially block other processes from making further progress.

Blocking algorithms introduce many problems such as:

- **Deadlock:** There is a circular lock-wait dependency among the processes, effectively preventing any processes from making progress.

- Convoy effect: One long process holding the lock will block other shorter processes contending for the lock.
- Priority inversion: A higher-priority process effectively has very low priority because it has to wait for another low priority process.

Furthermore, if a process that holds the lock dies, this will render the whole program unable to make any progress. This consideration holds even more weight in distributed computing because of a lot more failure modes, such as network failures, node failures, etc.

Therefore, while blocking algorithms, especially those using locks, are easy to write, they do not provide **progress guarantee** because **deadlock** or **livelock** can occur and their use of mutual exclusion is unnecessarily restrictive. Fortunately, there are other classes of algorithms which offer stronger progress guarantees.

2.4.2 Non-blocking algorithms

An algorithm is said to be **non-blocking** if a failure or slowdown in one process cannot cause the failure or slowdown in another process. Lock-free and wait-free algorithms are two especially interesting subclasses of non-blocking algorithms. Unlike blocking algorithms, they provide stronger degrees of progress guarantees.

2.4.2.1 Lock-free algorithms

Lock-free algorithms provide the following guarantee: Even if some processes are suspended, the remaining processes are ensured to make global progress and complete in bounded time. In other words, a process cannot cause hindrance to the global progress of the program. This property is invaluable in distributed computing; one dead or suspended process will not block the whole program, providing fault tolerance. Designing lock-free algorithms requires careful use of atomic instructions, such as Fetch-and-add (FAA), Compare-and-swap (CAS), etc which will be explained in Section 2.5.

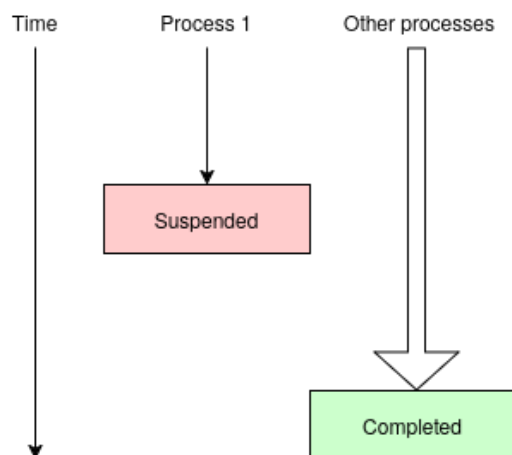


Figure 6: Lock-free algorithm: All the live processes together always finish in a finite amount of steps.

2.4.2.2 Wait-free algorithms

Wait-freedom offers the strongest degree of progress guarantee. It mandates that no process can cause constant hindrance to any running process. While lock-freedom ensures that at least one of the alive processes will make progress, wait-freedom guarantees that any alive process will finish in a finite number of steps. Wait-freedom can be desirable because it prevents starvation. Lock-freedom still allows the possibility of one process having to wait for another indefinitely, as long as some still make progress.

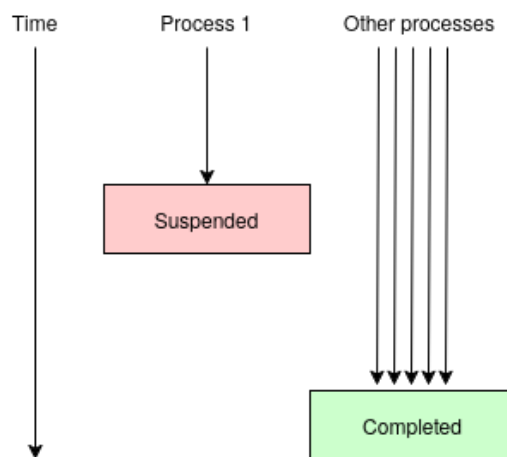


Figure 7: Wait-free algorithm: Any live process always finishes in a finite amount of steps.

2.5 Popular atomic instructions in designing non-blocking algorithms

As we have discussed in Section 2.4, blocking algorithms are not fault tolerant while non-blocking ones are, specifically lock-free and wait-free algorithms. Therefore, our design goal can be refined to linearizable non-blocking distributed MPSC queue. Techniques to achieve this is discussed next in this section. Issues, however, arise during the application of these techniques, whose resolution will be deferred to Section 2.6.

In non-blocking algorithms, finer-grained synchronization primitives than simple locks are required, which manifest themselves as atomic instructions. Therefore, it is necessary to get familiar with the semantics of these atomic instructions and common programming patterns associated with them.

2.5.1 Fetch-and-add (FAA)

Fetch-and-add (FAA) is a simple atomic instruction with the following semantics: It atomically increments a value at a memory location x by a and returns the previous value just before the increment. Informally, FAA's effect is equivalent to the function in Procedure 1, assuming that the function is executed atomically.

Procedure 1: `int fetch_and_add(int* x, int a)`

```
1 old_value = *x
2 *x = *x + a
3 return old_value
```

Fetch-and-add can be used to create simple distributed counters.

2.5.2 Compare-and-swap (CAS)

Compare-and-swap (CAS) is probably the most popular atomic operation instruction. The reason for its popularity is (1) CAS is a **universal atomic instruction** with the **consensus number** of ∞ , which means it is the most powerful atomic instruction [22] (2) CAS is implemented in most hardware (3) some concurrent lock-free data structures such as MPSC queues are more easily expressed using a powerful atomic instruction such as CAS.

The semantics of CAS is as follows. Given the instruction `CAS(memory location, old value, new value)`, atomically compares the value at `memory location` to see if it equals `old value`; if so, sets the value at `memory location` to `new value` and returns `true`; otherwise, leaves the value at `memory location` unchanged and returns `false`. Informally, its effect is equivalent to the function in Procedure 2.

Procedure 2: `bool compare_and_swap(int* x, int old_val, int new_val)`

```
1 if (*x == old_val)
2 |   *x = new_val
3 |   return true
4 return false
```

Compare-and-swap is very powerful and consequently, pervasive in concurrent algorithms and data structures.

Non-blocking concurrent algorithms often utilize CAS as follows. The steps 1-3 are retried until success.

1. Read the current value `old_value = read(memory location)`.
2. Compute `new_value` from `old_value` by manipulating some resources associated with `old_value` and allocating new resources for `new_value`.
3. Call `CAS(memory location, old_value, new_value)`. If that succeeds, the new resources for `new_value` remain valid because it was computed using valid resources associated with `old_value`, which has not been modified since the last read. Otherwise, free up the resources we have allocated for `new_value` because `old_value` is no longer there, so its associated resources are not valid.

This scheme is, however, susceptible to the ABA problem, which will be discussed in Section 2.6.1.

2.5.3 Load-link/Store-conditional (LL/SC)

Load-link/Store-conditional is actually a pair of atomic instructions for synchronization.

Semantically, load-link returns a value currently located at a memory location x while store-conditional sets the memory location x to a value v if there is no other write to x since the last load-link call, otherwise, the store-conditional call would fail.

Intuitively, LL/SC provides an easier synchronization primitive than CAS: LL/SC ensures that a store-conditional can only succeed if there is no access to a memory location, while CAS can still succeed in this case if the value at the memory location does not change. Due to this property, LL/SC is not vulnerable to the ABA problem (see Section 2.6.1). However, CAS is in fact as powerful as LL/SC, considering that they can implement each other [22].

Practically, store-conditional can still fail even if there is no write to the same memory location since the last load-link call. This is called a spurious failure. For example, consider the following generic sequence of events:

1. Thread X calls load-link on x and loads out v .
2. Thread X computes a new value v' .
3. Some *exceptional event* happens (discussed below). Assume that no other threads access x during this time.
4. Thread X calls store-conditional to store v' to x . It *should succeed* but *fails* anyway.

Exceptional events that can cause the store-conditional to fail spuriously include:

- Cache line flushing: If the cache line that caches the memory location x is written back to memory, logically, the memory location x has been accessed and therefore, the store-conditional fails.
- Context switch: If thread X is swapped out by the OS, cache lines may be invalidated and flushed out, which consequently leads to the first scenario.

LL/SC even though as powerful as CAS, is not as widespread as CAS; in fact, as of MPI-3, only CAS is supported.

2.6 Common issues when designing non-blocking algorithms

Atomic instructions are the option we choose when it comes to designing non-blocking algorithms (Section 2.5). However, there are problems usually associated with this approach, that is ABA problem (Section 2.6.1) and safe memory reclamation problem (Section 2.6.2). Proper solutions to these issues are required to complete our design process, which has been discussed at length in Section 2.2, Section 2.3, Section 2.4, Section 2.5. We move on to implementation techniques in section Section 2.7, Section 2.8, [15].

2.6.1 ABA problem

The ABA problem [23] is a notorious problem associated with the compare-and-swap atomic instruction. Because CAS is so widely used in non-blocking algorithms, the ABA problem almost has to always be accounted for.

As a reminder, here's how CAS is often utilized in non-blocking concurrent algorithms: The steps 1-3 are retried until success.

1. Read the current value `old value = read(memory location)`.
2. Compute new value from `old value` by manipulating some resources associated with `old value` and allocating new resources for new value.
3. Call `CAS(memory location, old value, new value)`. If that succeeds, the new resources for new value remain valid because it was computed using valid resources associated with `old value`, which has not been modified since the last read. Otherwise, free up the resources we have allocated for new value because `old value` is no longer there, so its associated resources are not valid.

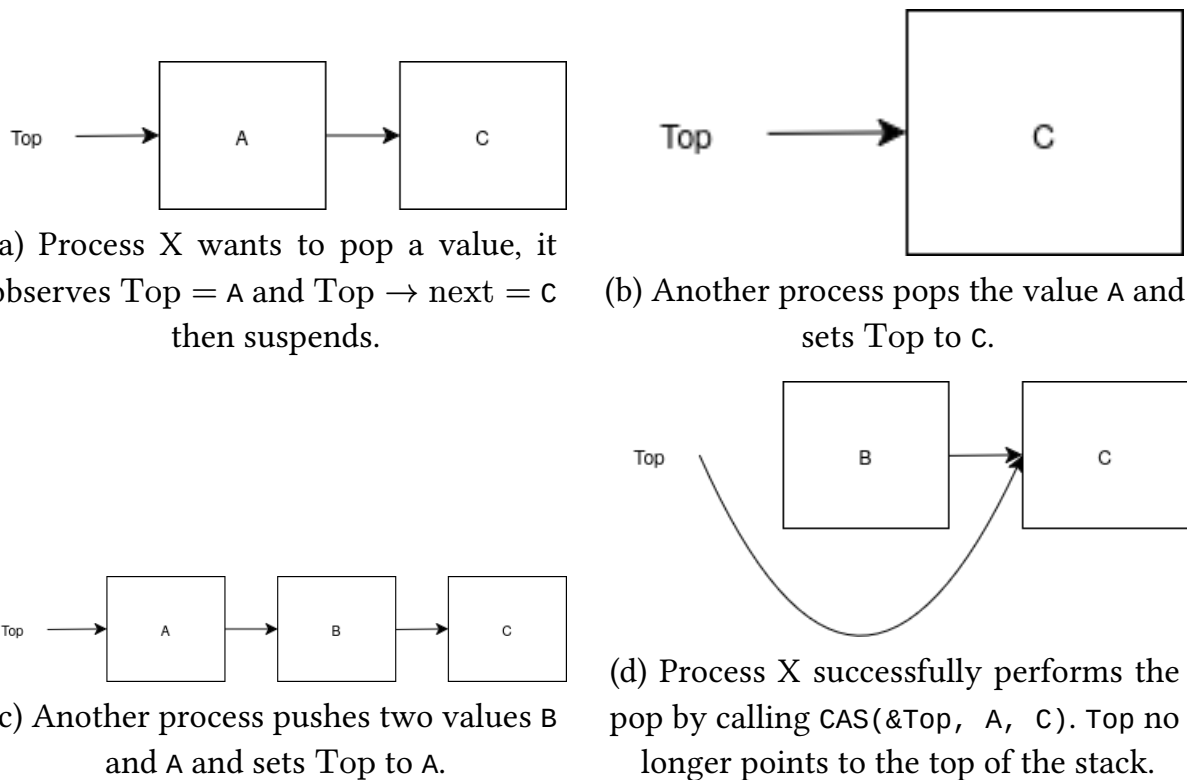


Figure 8: ABA problem in a linked-list stack.

As hinted, this scheme is susceptible to the notorious ABA problem. The following scenario illustrates an example of the ABA problem:

1. Process 1 reads the current value of memory location and reads out A.
2. Process 1 manipulates resources associated with A, and allocates resources based on these resources.
3. Process 1 suspends.
4. Process 2 reads the current value of memory location and reads out A.
5. Process 2 `CAS(memory location, A, B)` so that resources associated with A are no longer valid.

6. Process 3 CAS(memory location, B, A) and allocates new resources associated with A.
7. Process 1 continues and CAS(memory location, A, new value) relying on the fact that the old resources associated with A are still valid while in fact they aren't.

The ABA problem arises fundamentally because most algorithms assume a memory location is not accessed if its value is unchanged.

A specific case of the ABA problem is given in Figure 8.

To safeguard against the ABA problem, one must ensure that between the time a process reads out a value from a shared memory location and the time it calls CAS on that location, there is no possibility another process has CAS-ed the memory location to the same value.

A simple scheme that is widely used practically and also in this thesis is the **unique timestamp** scheme. This scheme's idea is simple: for each shared memory location that is affected by CAS operations, we reserve some bits of this memory location for a monotonic counter. Each time a CAS operation is carried out, this counter is incremented. Theoretically, the ABA problem would never happen because combining with this counter, the value of this memory location is always unique, due to the counter never repeating itself. However, practically, the counter can overflow and wrap around to the same value and the ABA problem would happen in this case. Therefore, the counter's range must be big enough so that this scenario can't virtually happen. Empirically, a counter of 32-bit should be enough. The drawback of this approach is that we have wasted 32 meaningful bits to avoid the ABA problem.

2.6.2 Safe memory reclamation problem

The problem of safe memory reclamation [24] often arises in concurrent algorithms that dynamically allocate memory. In such algorithms, dynamically-allocated memory must be freed at some point. However, there is a good chance that while a process is freeing memory, other processes contending for the same memory are keeping a reference to that memory. Therefore, deallocated memory can potentially be accessed, which is erroneous.

An example of unsafe memory reclamation is given in Figure 9.

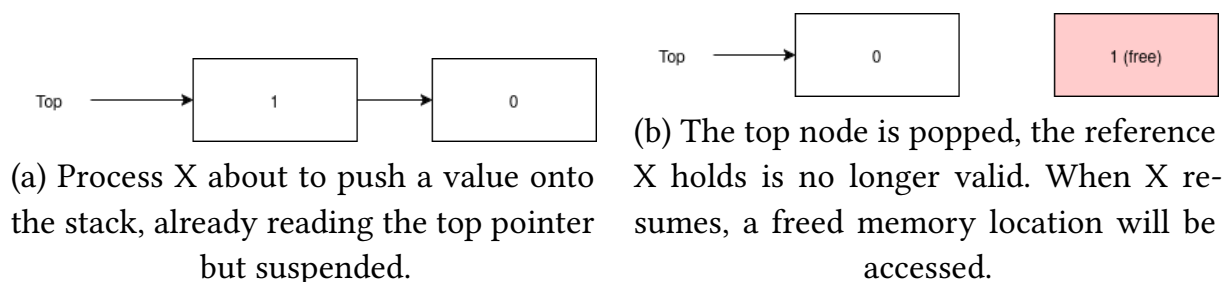


Figure 9: Unsafe memory reclamation in a LIFO stack.

Solutions to this problem must ensure that memory is only freed when no other processes are holding references to it. In garbage-collected programming environ-

ments, this problem can be conveniently pushed to the garbage collector. In non-garbage-collected programming environments, however, custom schemes must be utilized.

2.7 MPI-3 - A popular distributed programming library interface specification

To implement the design linearizable non-blocking distributed MPSC queues, the most basic choice we can go with is MPI, which will be the matter of this section. We specifically focus on the MPI-3 RMA API, because as will be explained, it facilitates the easy implementation of irregular applications such as MPSC queues. An approach of using MPI-3 RMA is covered in Section 2.8. More advanced implementation techniques will be covered in Section 2.9.

MPI stands for message passing interface, which is a **message-passing library interface specification**. Design goals of MPI include high availability across platforms, efficient communication, thread safety, reliable and convenient communication interface while still allowing hardware-specific accelerated mechanisms to be exploited [2].

2.7.1 MPI-3 RMA

RMA in MPI RMA stands for remote memory access. RMA APIs were introduced in MPI-2 and their capabilities are further extended in MPI-3 to conveniently express irregular applications [20]. In general, RMA is intended to support applications with dynamically changing data access patterns where the data distribution is fixed or slowly changing [2]. This is very similar to the properties of irregular applications as discussed in Section 2.1. In such applications, one process, based on the data it needs, knowing the data distribution, can compute the nodes where the data is stored. However, because the data access pattern is not known, each process cannot know whether any other processes will access its data. Using the traditional Send/Receive interface, both sides need to issue matching operations by distributing appropriate transfer parameters. This is not suitable, as previously explained; only the side that needs to access the data knows all the transfer parameters while the side that stores the data cannot anticipate this.

2.7.2 MPI-RMA communication operations

RMA only requires one side to specify all the transfer parameters and thus only that side to participate in data communication [2].

To utilize MPI RMA, each process needs to open a memory window to expose a segment of its memory to RMA communication operations such as remote writes (MPI_PUT), remote reads (MPI_GET) or remote accumulates (MPI_ACCUMULATE, MPI_GET_ACCUMULATE, MPI_FETCH_AND_OP, MPI_COMPARE_AND_SWAP) [2]. These remote communication operations only require one side to specify.

2.7.3 MPI-RMA synchronization

Besides communication of data from the sender to the receiver, one also needs to synchronize the sender with the receiver. That is, there must be a mechanism to ensure the completion of RMA communication calls or that any remote operations have taken effect. For this purpose, MPI RMA provides **active target synchronization** and **passive target synchronization**. In this document, we are particularly interested in **passive target synchronization** as this mode of synchronization does not require the target process of an RMA operation to explicitly issue a matching synchronization call with the origin process, easing the expression of irregular applications [20].

In **passive target synchronization**, any RMA communication calls must be within a pair of `MPI_Win_lock`/`MPI_Win_unlock` or `MPI_Win_lock_all`/`MPI_Win_unlock_all`. After the unlock call, those RMA communication calls are guaranteed to have taken effect. One can also force the completion of those RMA communication calls without the need for the call to unlock using flush calls such as `MPI_Win_flush` or `MPI_Win_flush_local`.



Figure 10: An illustration of passive target communication. Dashed arrows represent synchronization (source: [2]).

2.8 Pure MPI - A porting approach of shared memory algorithms to distributed algorithms

With MPI (Section 2.7), we have the most basic facility to adapt shared-memory algorithms to distributed algorithms, which is MPI-3 RMA. However, MPI-3 RMA offer a wide range of utilities, which may not be quite well-suited to implement non-blocking distributed MPSC queues. In this section, we introduce one technique utilizing MPI-3 RMA to implement distributed non-blocking distributed algorithms. The BCL CoreX library (Section 2.9) is built on top of this approach.

In pure MPI, we use MPI exclusively for communication and synchronization. With MPI RMA, the communication calls that we utilize are [2]:

- Remote read: MPI_Get
- Remote write: MPI_Put
- Remote accumulation: MPI_Accumulate, MPI_Get_accumulate, MPI_Fetch_and_op and MPI_Compare_and_swap.

For lock-free synchronization, we choose to use **passive target synchronization** with MPI_Win_lock_all/MPI_Win_unlock_all.

In the MPI-3 specification [2], these functions are specified as in Table 1.

Operation	Usage
MPI_Win_lock_all	Starts an RMA access epoch to all processes in a memory window, with a lock type of MPI_LOCK_SHARED. The calling process can access the window memory on all processes in the memory window using RMA operations. This routine is not collective.
MPI_Win_unlock_all	Matches with an MPI_Win_lock_all to unlock a window previously locked by that MPI_Win_lock_all.

Table 1: Specification of MPI_Win_lock_all and MPI_Win_unlock_all.

The reason we choose this is 3-fold:

- Unlike **active target synchronization**, **passive target synchronization** does not require the process whose memory is being accessed by an MPI RMA communication call to participate. This is in line with our intention to use MPI RMA to easily model irregular applications like MPSC queues.
- Unlike **active target synchronization**, MPI_Win_lock_all and MPI_Win_unlock_all do not need to wait for a matching synchronization call in the target process, and thus, are not delayed by the target process.
- Unlike **passive target synchronization** with MPI_Win_lock/MPI_Win_unlock, multiple calls of MPI_Win_lock_all can succeed concurrently, so one process needing to issue MPI RMA communication calls does not block others.

An example of our pure MPI approach with MPI_Win_lock_all/MPI_Win_unlock_all, inspired by [20], is illustrated in the following:

```

MPI_Win_lock_all(0, win);

MPI_Get(...); // Remote get
MPI_Put(...); // Remote put
MPI_Accumulate(..., MPI_REPLACE, ...); // Atomic put
MPI_Get_accumulate(..., MPI_NO_OP, ...); // Atomic get
MPI_Fetch_and_op(...); // Remote fetch-and-op
MPI_Compare_and_swap(...); // Remote compare and swap
...

MPI_Win_flush(...); // Make previous RMA operations take effect
MPI_Win_flush_local(...); // Make previous RMA operations take effect
locally
...

MPI_Win_unlock_all(win);

```

Listing 3: An example snippet showcasing our synchronization approach in MPI RMA.

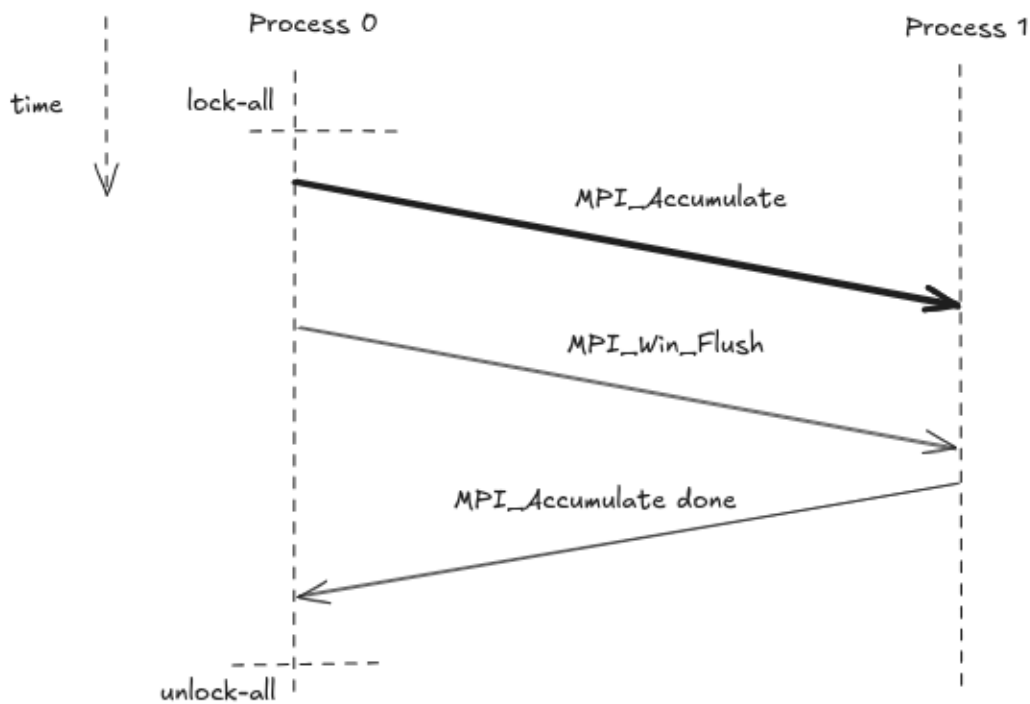


Figure 11: An illustration of our synchronization approach in MPI RMA.

2.9 BCL CoreX

BCL CoreX [15] is a high-level library built on top of MPI to facilitate the design of non-blocking algorithms for distributed-memory machines. In principle, it utilizes the pure MPI approach that we have covered in Section 2.8.

A subset of the primitives provided by BCL CoreX is presented below. We will utilize these primitives in our algorithm specification.

gptr<T>

A global pointer that points to a variable of type T . A global pointer is one that can point to a variable outside of the current process's address space. The process whose address space a global pointer points to is called the *host*. Like normal pointers, pointer arithmetic also works with global pointers, which allows global pointers to point to remote arrays.

T read(gptr<T> ptr, T* dest)

Issue a synchronous read on the location pointed to by `ptr` and stores the read value in `dest`.

T write(gptr<T> ptr, T* src)

Issue a synchronous write on the location pointed to by `ptr` that writes the value stored in `src`.

T faa(gptr<T> ptr, T inc)

Issue a synchronous fetch-and-add operation on the location pointed to by `ptr`.
 T must be an integral type of less than 64 bits.

T cas(gptr<T> ptr, T old_val, T new_val)

Issue a synchronous compare-and-swap operation on the location pointed to by `ptr`.
 T must be a type of less than 64 bits.

A remote operation occurs when one of the primitive operations is applied on the global pointer that points to a non-local address space. Otherwise, a local operation occurs. Typically, remote operations are very expensive compared to local operations.

Chapter III Related works

Given our decision to adapt shared-memory data structures for creating non-blocking distributed MPSC queues with BCL CoreX (Section 2.9), we next explore existing non-blocking shared-memory MPSC queues in Section 3.1. Additionally, we analyze current distributed MPSC queue implementations in Section 3.2 to establish benchmarking baseline. Based on this analysis, we selected LTQueue as our candidate shared-memory MPSC queue for distributed adaptation and chose AMQueue as our primary benchmark reference.

3.1 Non-blocking shared-memory MPSC queues

There exists numerous research into the design of non-blocking shared memory MPSCs. Interestingly, research into non-blocking MPSC queues is noticeably scarce. In reality, we have only found 4 papers that are concerned with the direct support of non-blocking MPSC queues: LTQueue [12], DQueue [11], WRLQueue [10], and Jiffy [13]. Table 2 summarizes the characteristics of these algorithms.

MPSC queues	LTQueue [12]	DQueue [11]	WRLQueue [10]	Jiffy [13]
ABA solution	Load-link/ Store-conditional	Incorrect custom scheme (*)	Custom scheme	Custom scheme
Safe memory reclamation	Custom scheme	Incorrect custom scheme (*)	Custom scheme	Insufficient custom scheme
Progress guarantee of dequeue	Wait-free	Wait-free	Blocking (*)	Wait-free
Progress guarantee of enqueue	Wait-free	Wait-free	Wait-free	Wait-free

Table 2: Summary of existing shared memory MPSC queues. The cell marked with (*) indicates that our evaluation contradicts with the authors' claims.

3.1.1 LTQueue

To our knowledge, LTQueue [12] is the earliest paper that directly focuses on the design of a wait-free shared memory MPSC queue.

This algorithm is wait-free with $O(\log n)$ time complexity for both enqueues and dequeues, with n being the number of enqueueers due to a novel timestamp-update scheme and a tree-structure organization of timestamps.

The basic structure of LTQueue is given in Figure 12. In LTQueue, each enqueueer maintains an SPSC queue that only it and the dequeuer can access. This SPSC queue must additionally support the readFront operation, which returns the front element

currently in the SPSC. The SPSC can be any implementation that conforms to this interface. In the original paper, the SPSC is represented as a simple linked list.

The rectangular nodes at the bottom in Figure 12 represent an enqueueer, whose SPSC contains items with 2 fields: value and timestamp. Every enqueueer has to timestamp its data before enqueueing. The timestamps can be obtained using a distributed counter shared by all the enqueueers.

The purpose of timestamping is to determine the order to dequeue the items from the local SPSCs. To efficiently maintain the timestamps and determine which SPSC to dequeue from first, a tree structure with a min-heap property is built upon the enqueueer nodes. The original algorithm leaves the exact representation of the tree open, for example, the arity of the tree, which is shown to be 2 in Figure 12. The circle-shaped nodes in this figure represent the nodes in this tree structure, which are shared by all processes. Each node stores the minimum timestamp along with the owner enqueueer's rank (an identifier given to a process) in the subtree rooted at that node. After every modification to the local SPSC, i.e., an enqueue or a dequeue, the changes must be propagated up to the root node.

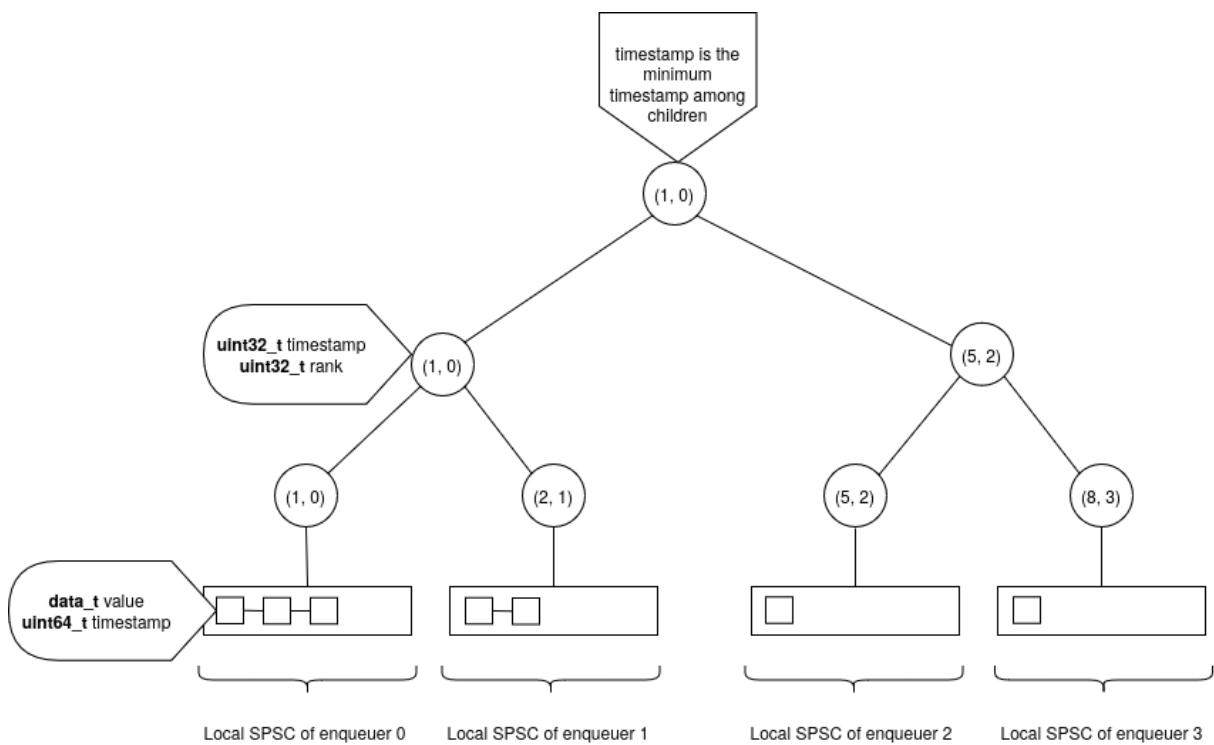


Figure 12: LTQueue's structure.

To dequeue, the dequeuer simply looks at the root node to determine the rank of the enqueueer to dequeue its SPSC.

The fundamental idea contributing to LTQueue's wait-freedom is the wait-free timestamp-propagation procedure. If there is a change to an enqueueer's SPSC, the timestamp of any nodes that lie on the path from the enqueueer to the root node is refreshed. The timestamp-refreshing procedure is simple:

- Call load-link on the node's (timestamp, rank).

- Look at all the timestamps of the node's children and determine the minimum timestamp and its owner rank.
- Call store-conditional to store the new minimum timestamp and the new owner rank to the current node.

Notice that due to contention, the timestamp-refreshing procedure can fail. In that case, the timestamp-propagation procedure simply retries the timestamp-refreshing procedure one more time. This second call, again, can fail. However, after this second call, the node's timestamp is guaranteed to be up-to-date. The intuition behind this is demonstrated in Figure 13. Furthermore, because every node is refreshed at most twice, the timestamp-refresh procedure should finish in a finite number of steps.



Figure 13: Intuition on how timestamp-refreshing works.

The LTQueue algorithm avoids ABA entirely by utilizing load-link/store-conditional. This represents a challenge to directly implementing this algorithm in a distributed environment.

The memory reclamation responsibility is handled by the SPSC structure, which is fairly trivial with a custom scheme.

The design of each enqueueer maintaining a separate SPSC allows multiple enqueueers to successfully enqueue their data in parallel without stepping on each other's toes. This can potentially scale well to a large number of processes. However, scalability may be limited due to potentially growing contention during timestamp propagation. The performance of LTQueue in shared-memory environments may still have a lot of room for improvement, i.e., more cache-aware design, avoiding unnecessary contention, etc. Nevertheless, its timestamp-refreshing scheme is interesting in and of itself and can potentially inspire the design of new algorithms. In fact, LTQueue's idea is core to one of our optimized distributed MPSC queue algorithms, Slotqueue (Section 4.3).

3.1.2 DQueue

DQueue [11] focuses on optimizing performance, aiming to be cache-friendly and avoid expensive atomic instructions such as CAS.

The basic structure of DQueue is demonstrated in Figure 14.

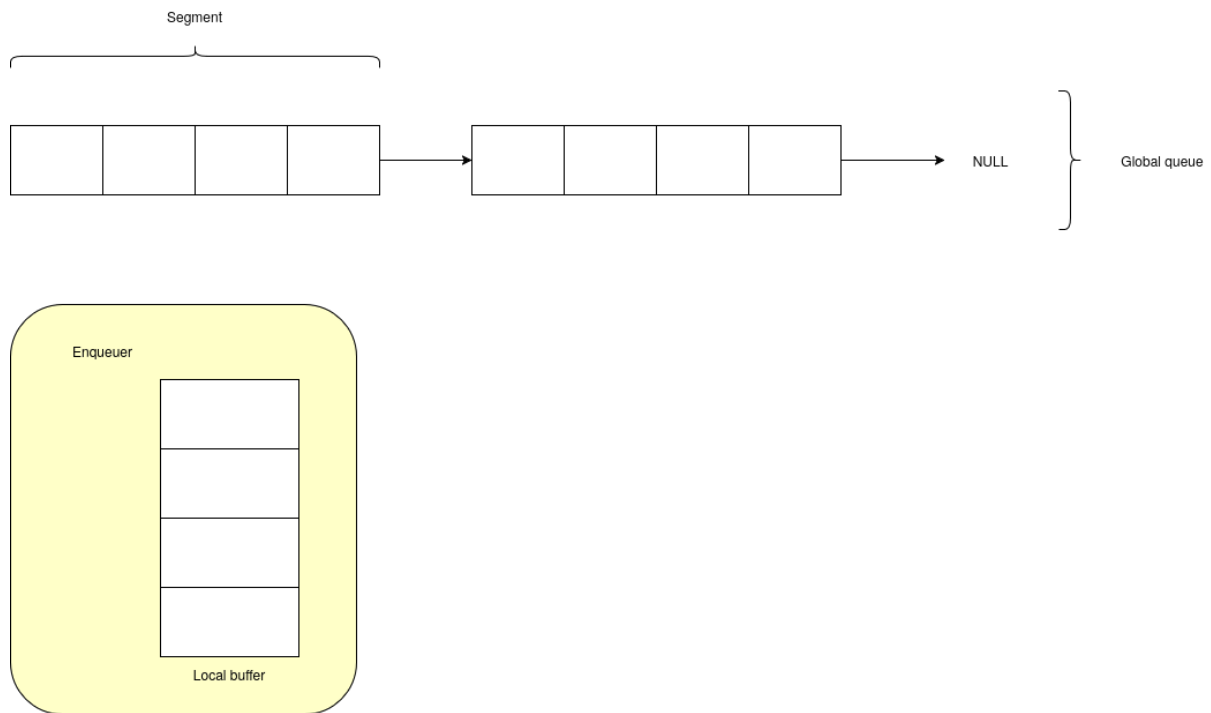


Figure 14: DQueue's structure.

The global queue where data is stored is represented as a linked list of segments. A segment is simply a contiguous array of data items. This design allows for unbounded queue capacity while still allowing a fair degree of random access within a segment. This allows us to use indices to index elements in the queue, thus permitting the use of inexpensive FAA instructions to swing the head and tail indices.

Each enqueueer maintains a local buffer to batch enqueue items into before flushing to the global queue. This helps prevent contention and plays nicely with the cache. To enqueue an item, an enqueueer simply FAA-s the head index to reserve a slot in the global queue; the obtained index is stored along with the data in the local buffer so that when flushing the local buffer, the enqueueer knows where to write the data into the global queue. Note that while flushing, an index may point to a not-yet-existent slot in the global queue. Therefore, new segments must be allocated on the fly and CAS-ed to the end of the queue.

The dequeuer dequeues the items by looking at the head index. If the queue is not empty but the slot at the head index is empty, the dequeuer utilizes a helping mechanism by looking at all enqueueers to help them flush out their local buffer. After this, the head slot is guaranteed to be non-empty, and the dequeuer can finally dequeue this value.

The ABA problem is solved by relying on its safe memory reclamation scheme. In DQueue, CAS is only used to update the tail pointer to point to a newly allocated segment. Therefore, the ABA problem in DQueue only involves internal manipulation

of pointers to dynamically allocated memory. This means that if a proper memory reclamation scheme is used, the ABA problem cannot occur.

DQueue relies on a dedicated garbage collection thread to reclaim segments that have been exhausted by the dequeuer. However, this should be a careful process as even though some segments have been exhausted, some enqueueers can still hold an index that references one of these segments. DQueue implements this by reclaiming all exhausted segments if there is no enqueueer holding an index referencing these segments. Unfortunately, we believe DQueue's scheme is unsafe. Specifically, as described, DQueue allows the garbage collection thread to reclaim non-adjacent segments in the global queue without patching any of the next pointers. Any segment just before a reclaimed segment would point to a deallocated next segment. By definition, this segment was not reclaimed because it is referenced by an enqueueer. This means this enqueueer cannot traverse the next pointer chain to get to the end of the queue without accessing an already-deallocated segment.

If adapted to a distributed environment, the flushing may be expensive, both from the point of view of the enqueueer and the dequeuer. If the dequeuer has to help every enqueueer to flush their local buffer, which should always result in at least one remote operation, the cost would be prohibitively high. Similarly, each flush requires the enqueueer to issue at least one remote operation, but this is at least acceptable as flushing is infrequent.

Still, we can see that the pattern of maintaining a local buffer inside each enqueueer repeats throughout the literature, which we can definitely apply when designing distributed MPSC queues.

3.1.3 WRLQueue

WRLQueue [10] is a lock-free MPSC queue specifically designed for embedded real-time systems. Its main purpose is to avoid excessive modification of storage space.

WRLQueue is simply a pair of buffers: one is worked on by multiple enqueueers, and the other is worked on by the dequeuer. The structure of WRLQueue is shown in Figure 15.



Figure 15: WRLQueue's structure.

The enqueueers batch their enqueues and write multiple elements onto the buffer at once. They use the usual scheme of FAA-ing the tail index (`write_pos` in Figure 15) to reserve their slots and write data items at their leisure.

The dequeuer, upon invocation, will swap its buffer with the enqueueers' buffer to dequeue from it, as in Figure 16. However, WRLQueue explicitly states that the dequeuer has to wait for all enqueue operations to complete in the other buffer before swapping. If an enqueue suspends or dies, the dequeuer will experience a slowdown; this clearly violates the property of non-blocking. Therefore, we believe that WRLQueue is blocking, concerning its dequeue operation.

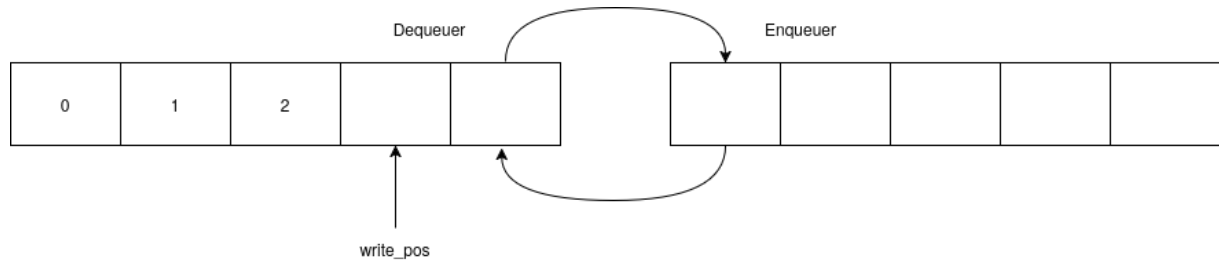


Figure 16: WRLQueue's dequeue operation

3.1.4 Jiffy

Jiffy [13] is a fast and memory-efficient wait-free MPSC queue by avoiding excessive allocation of memory.

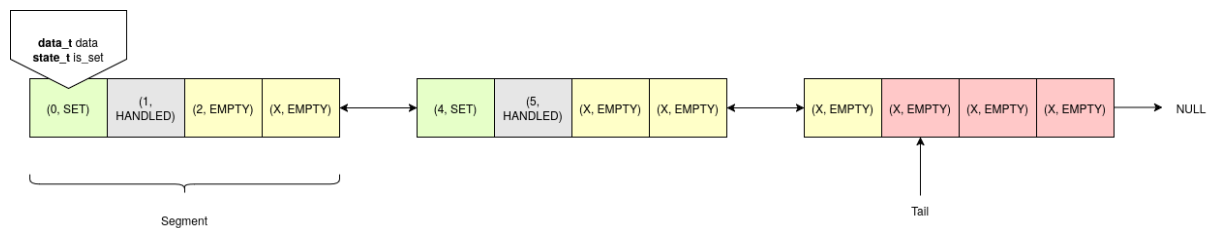


Figure 17: Jiffy's structure.

Like DQueue, Jiffy represents the queue as a doubly-linked list of segments as in Figure 17. This design again allows Jiffy to be unbounded while using head and tail indices to index elements. Each segment contains a pointer to a dynamically allocated array of slots, instead of directly storing the array. Each slot in the segment contains the data item and a state of that slot (`state_t` in the figure). There are 3 states: SET, EMPTY, and HANDLED. Initially, all slots are EMPTY. Instead of keeping a global head index, there are per-segment Head indices pointing to the first non-HANDLED slot. However, there is still one global Tail index shared by all the processes.

To enqueue, each enqueueer would FAA the Tail to reserve a slot. If the slot isn't in the linked list yet, it tries to allocate new segments and CAS them at the end of the linked list until the slot is available. It then traverses to the desired segment by following the previous pointers starting from the last segment. It then writes the data and sets the slot's state to SET. Notice that EMPTY slots actually have two substates. If an EMPTY slot is before the Tail index, that slot is actually reserved by an enqueueer but has not been set yet, while the EMPTY slots after the Tail index are truly empty.

To dequeue, the dequeuer would start from the Head index of the first segment, scanning until it finds the first non-HANDLED slot before the end of the queue. If there is

no such slot, the queue is empty, and the dequeuer would return nothing. If this slot is SET, it simply reads the data item in this slot and sets it to HANDLED. If this slot is EMPTY, that means this slot has been reserved by an enqueueer that has not finished. In this case, the dequeuer performs a scan forward to find the first SET slot. If not found, the dequeuer returns nothing. Otherwise, it continues to repeatedly scan all slots between the first non-HANDLED and the last found SET slot until the first SET slot in this interval is unchanged between 2 scans. Only then, the dequeuer would return the data item in this SET slot and mark it as HANDLED.

Similar to DQueue, CAS is only used when appending new segments at the end of the queue. Therefore, the ABA problem only involves internal manipulation of pointers to dynamically allocated memory. Consequently, if a proper memory reclamation scheme is utilized, the ABA problem is also properly solved.

Regarding memory reclamation, Jiffy does not specify a sufficient scheme: If one enqueueer is delayed forever, no memory is ever reclaimed. As a consequence, if an enqueueer is delayed for too long, the system will run out of memory, causing other enqueueers to fail without making any progress. Effectively, Jiffy is not wait-free.

3.1.5 Remarks

Out of the 4 investigated MPSC queue algorithms, we quickly eliminate DQueue, WRLQueue, and Jiffy as potential candidates for porting to a distributed environment because they either do not provide a sufficient progress guarantee or protection against the ABA problem and memory reclamation problem. Therefore, we will only adapt LTQueue for distributed environments in the next section. LTQueue also presents some challenges, though, as it utilizes LL/SC for the ABA solution, which does not exist in distributed environments. Consequently, to adapt LTQueue, we have to work around LTQueue's usage of LL/SC.

3.2 Distributed MPSC queues

This section summarizes, to the best of our knowledge, existing MPSC queue algorithms, which is reflected in Section 3.2.

The only paper we have found so far that either mentions directly or indirectly the design of an MPSC queue is [1]. [1] introduces a hosted blocking (the original paper claims that it is lock-free) bounded distributed MPSC queue called active-message queue (AMQueue) that bears resemblance to WRLQueue in [10].

FIFO queues	Active-message queue (AMQueue) [1]
Progress guarantee of dequeue	Blocking (*)
Progress guarantee of enqueue	Wait-free
ABA solution	No compare-and-swap usage
Safe memory reclamation	Custom scheme

Table 3: Characteristic summary of existing distributed MPSC queues.

R stands for remote operations and L stands for local operations.

(*) [1] claims that it is lock-free.

The structure of AMQueue is given in Figure 18. The MPSC is split into 2 queues, each maintaining its own set of control variables:

- **WriterCnt**: The number of enqueueers currently writing in this queue.
- **Offset**: The index to the first empty entry in the queue. Note that any shared data and control variables are hosted on the dequeuer.

To determine which queue to read and write, the **QueueNum** binary variable is used. If **QueueNum** is 0, then the first queue is being actively written by enqueueers, and the second queue is being reserved for the dequeuer, and vice versa.

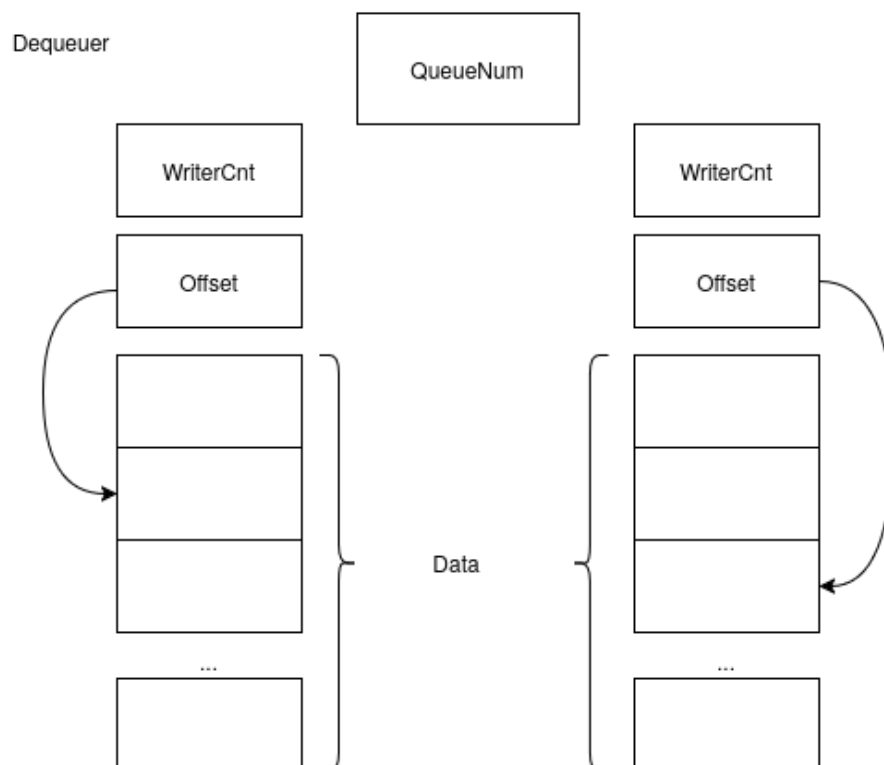


Figure 18: AMQueue's structure.

To enqueue, the enqueueer first reads the **QueueNum** variable to see which queue is active. The enqueueer then registers for that queue by atomically FAA-ing the corresponding **WriterCnt** variable. If the fetched value is negative, though, the **QueueNum** queue is

being swapped for dequeuing, and the enqueueer has to decrement the `WriterCnt` variable and repeat the process until `WriterCnt` is positive. After a successful registration, the enqueueer then reserves an entry in the data array by FAA-ing the `Offset` variable. After that, the enqueueer can enqueue data at its leisure. Upon success, the enqueueer has to decrement `WriterCnt` before returning.

To dequeue, the dequeuer inverts `QueueNum` to direct future enqueueers to the other queue. The dequeuer then subtracts a sufficiently large number from `WriterCnt` to signal to other enqueueers that it has started processing. The dequeuer has to wait for all current enqueueers in the queue to finish by repeatedly checking the `WriterCnt` variable, hence the blocking property. After all enqueueers have finished, the dequeuer then batch-dequeues all data in the queue, resetting the `Offset` and `WriterCnt` variables to 0.

Based on our discussion, there is currently no non-blocking distributed MPSC queue in the literature. This makes our research the first one of its kind to be about non-blocking distributed MPSC queues. `AMQueue` will serve as a benchmarking baseline for our MPSC queues in Chapter V.

Chapter IV Distributed MPSC queues

Based on the MPSC queue algorithms we have surveyed in Chapter III, we propose two wait-free distributed MPSC queue algorithms:

- dLTQueue (Section 4.2) is a direct modification of the original LTQueue [12] without any usage of LL/SC, adapted for distributed environment.
- Slotqueue (Section 4.3) is inspired by the timestamp-refreshing idea of dLTQueue [12] and repeated-rescan of Jiffy [13]. Although it still bears some resemblance to LTQueue, we believe that it is more optimized for distributed context.

In actuality, dLTQueue and Slotqueue are more than simple MPSC algorithms. They are *MPSC queue wrappers*, that is, given an SPSC queue implementation, they yield an MPSC implementation. There is one additional constraint: The SPSC interface must support an additional readFront operation, which returns the first data item currently in the SPSC queue.

This fact has an important implication: when we are talking about the characteristics (correctness, progress guarantee, performance model, ABA solution and safe memory reclamation scheme) of an MPSC queue wrapper, we are talking about the correctness, progress guarantee, performance model, ABA solution and safe memory reclamation scheme of the wrapper that turns an SPSC queue to an MPSC queue:

- If the underlying SPSC queue is linearizable, the resulting MPSC queue is linearizable.
- The resulting MPSC queue's progress guarantee is the weaker guarantee between the wrapper's and the underlying SPSC's.
- If the underlying SPSC queue is safe against ABA problem and memory reclamation, the resulting MPSC queue is also safe against these problems.
- If the underlying SPSC queue is unbounded, the resulting MPSC queue is also unbounded.
- The theoretical performance of dLTQueue and Slotqueue has to be coupled with the theoretical performance of the underlying SPSC.

The characteristics of these MPSC queue wrappers are summarized in Table 4. For benchmarking purposes, we use a baseline distributed SPSC introduced in Section 4.1 in combination with the MPSC queue wrappers. The characteristics of the resulting MPSC queues are also shown in Table 4, which will be proven in Section Appendix A.

MPSC queues	dLTQueue	Slotqueue
Correctness	Linearizable	Linearizable
Progress guarantee of dequeue	Wait-free	Wait-free
Progress guarantee of enqueue	Wait-free	Wait-free
Dequeue time-complexity (*)	$4 \log_2(n)R + 6 \log_2(n)L$	$3R + 2nL$
Enqueue time-complexity (*)	$6 \log_2(n)R + 4 \log_2(n)L$	$4R + 3L$
ABA solution	Unique timestamp	No hazardous ABA problem
Safe memory reclamation	No dynamic memory allocation	No dynamic memory allocation

Table 4: Characteristic summary of our proposed distributed MPSC queues.

(1) n is the number of enqueueers.

(2) R stands for **remote operation** and L stands for **local operation**.

(*) The underlying SPSC is assumed to be our simple distributed SPSC in Section 4.1.

The rest of this chapter is organized as follows. Section 4.1 describes a simple baseline distributed SPSC that is utilized as the underlying SPSC in our MPSC queues. Section 4.2 and Section 4.3 introduce dLTQueue and Slotqueue, our two wait-free MPSC queues that are our main contributions in this thesis.

In these next few descriptions, we assume that each process in our program is assigned a unique number as an identifier, which is termed as its **rank**. The numbers are taken from the range of $[0, \text{size} - 1]$, with size being the number of processes in our program.

4.1 A simple baseline distributed SPSC

The two MPSC queue wrapper algorithms we propose in Section 4.2 and Section 4.3 both utilize a baseline distributed SPSC data structure, which we will present in this section.

For implementation simplicity, we present a bounded SPSC, effectively make our proposed algorithms support only a bounded number of elements. However, one can trivially substitute another distributed unbounded SPSC to make our proposed algorithms support an unbounded number of elements, as long as this SPSC supports the same interface as ours.

The SPSC queue uses a circular array `Data` with a fixed `Capacity`. It maintains two indices: `First` marks the oldest item not yet removed, and `Last` marks the next available slot for insertion. Both indices use modulo arithmetic (`First % Capacity` and `Last % Capacity`) to wrap around the array.

For performance optimization, each process maintains local cached copies of these indices in `First_buf` and `Last_buf`. All indices start at zero. The memory layout is distributed between processes: the dequeuer hosts the `First` and `Last` indices, while the enqueueer hosts the `Data` array itself.

Placement-wise, all queue data in this SPSC is hosted on the enqueueer while the control variables i.e. `First` and `Last`, are hosted on the dequeuer.

Shared variables

```
Data: gptr<data_t>
First: gptr<uint64_t>
Last: gptr<uint64_t>
```

Enqueueer-local variables

```
First_buf: uint64_t
Last_buf: uint64_t
Capacity: uint64_t
```

Dequeuer-local variables

```
First_buf: uint64_t
Last_buf: uint64_t
Capacity: uint64_t
```

The procedures of the enqueueer are given as follows.

Procedure 4: `bool spsc_enqueue(data_t v)`

```
1 new_last = Last_buf + 1
2 if (new_last - First_buf > Capacity)
3   read(First, &First_buf)
4   if (new_last - First_buf > Capacity)
5     | return false
6 write(Data + Last_buf % Capacity, &v)
7 write>Last, &new_last)
8 Last_buf = new_last
9 return true
```

`spsc_enqueue` first computes the new `Last` value (Line 1). If the queue is full as indicated by the difference the new `Last` value and `First_buf` (Line 2), there can still be the possibility that some elements have been dequeued but `First_buf` has not been synced with `First` yet, therefore, we first refresh the value of `First_buf` by fetching from `First` (Line 3). If the queue is still full (Line 4), we signal failure (Line 5). Otherwise, we proceed to write the enqueued value to the entry at `Last_buf % Capacity` (Line 6), increment `Last` (Line 7), update the value of `Last_buf` (Line 8) and signal success (Line 9).

Procedure 5: `bool spsc_readFronte(data_t* output)`

```
10 if (First_buf >= Last_buf)
11   | return false
12 read(First, &First_buf)
13 if (First_buf >= Last_buf)
14   | return false
15 read(Data + First_buf % Capacity, output)
16 return true
```

`spsc_readFronte` first checks if the SPSC is empty based on the difference between `First_buf` and `Last_buf` (Line 10). Note that if this check fails, we signal failure immediately (Line 11) without refetching either `First` or `Last`. This suffices because `Last` cannot be out-of-sync with `Last_buf` as we are the enqueueer and `First` can only increase since the last refresh of `First_buf`, therefore, if we refresh `First` and `Last`, the condition on Line 10 would return false anyways. If the SPSC is not empty, we refresh `First` and re-perform the empty check (Line 13 - Line 14). If the SPSC is again not empty, we read the queue entry at `First_buf % Capacity` into output (Line 15) and signal success (Line 16).

The procedures of the dequeuer are given as follows.

Procedure 6: `bool spsc_dequeue(data_t* output)`

```
17 new_first = First_buf + 1
18 if (new_first > Last_buf)
19   | read>Last, &Last_buf)
20   | if (new_first > Last_buf)
21     | return false
22 read(Data + First_buf % Capacity, output)
23 write(First, &new_first)
24 First_buf = new_first
25 return true
```

`spsc_dequeue` first computes the new `First` value (Line 17). If the queue is empty as indicated by the difference the new `First` value and `Last_buf` (Line 18), there can still be the possibility that some elements have been enqueued but `Last_buf` has not been synced with `Last` yet, therefore, we first refresh the value of `Last_buf` by fetching from `Last` (Line 19). If the queue is still empty (Line 20), we signal failure (Line 21). Otherwise, we proceed to read the top value at `First_buf % Capacity` (Line 22) into output, increment `First` (Line 23) - effectively dequeue the element, update the value of `First_buf` (Line 24) and signal success (Line 25).

Procedure 7: `bool spsc_readFrontd(data_t* output)`

```
26 if (First_buf >= Last_buf)
27   read(Last, &Last_buf)
28   if (First_buf >= Last_buf)
29     | return false
30 read(Data + First_buf % Capacity, output)
31 return true
```

`spsc_readFrontd` first checks if the SPSC is empty based on the difference between `First_buf` and `Last_buf` (Line 26). If this check fails, we refresh `Last_buf` (Line 27) and recheck (Line 28). If the recheck fails, signal failure (Line 29). If the SPSC is not empty, we read the queue entry at `First_buf % Capacity` into `output` (Line 30) and signal success (Line 31).

4.2 dLTQueue - Straightforward LTQueue adapted for distributed environment

This algorithm presents our most straightforward effort to port LTQueue [12] to distributed context. The main challenge is that LTQueue uses LL/SC as the universal atomic instruction and also an ABA solution, but LL/SC is not available in distributed programming environments. We have to replace any usage of LL/SC in the original LTQueue algorithm. We use compare-and-swap and the well-known monotonic time-stamp scheme to guard against ABA problem.

4.2.1 Overview

The structure of our dLTQueue is shown as in Figure 19.

We differentiate between 2 types of nodes: *enqueueer nodes* (represented as the rectangular boxes at the bottom of Figure 19) and normal *tree nodes* (represented as the circular boxes in Figure 19).

Each enqueueer node corresponds to an enqueueer. Each time the local SPSC is enqueued with a value, the enqueueer timestamps the value using a distributed counter shared by all enqueueers. An enqueueer node stores the SPSC local to the corresponding enqueueer and a `min_timestamp` value which is the minimum timestamp inside the local SPSC.

Each tree node stores the rank of an enqueueer process. This rank corresponds to the enqueueer node with the minimum timestamp among the node's children's ranks. The tree node that is attached to an enqueueer node is called a *leaf node*, otherwise, it is called an *internal node*.

Note that if a local SPSC is empty, the `min_timestamp` variable of the corresponding enqueueer node is set to `MAX_TIMESTAMP` and the corresponding leaf node's rank is set to `DUMMY_RANK`.



Figure 19: dLTQueue's structure.

Placement-wise:

- The enqueueer nodes are hosted at the corresponding enqueueer.
- All the tree nodes are hosted at the dequeueer.
- The distributed counter, which the enqueueers use to timestamp their enqueued value, is hosted at the dequeueer.

4.2.2 Data structure

Below is the types utilized in dLTQueue.

Types

`data_t` = The type of the data to be stored.

`spsc_t` = The type of the SPSC, this is assumed to be the distributed SPSC in Section 4.1.

`rank_t` = The type of the rank of an enqueueer process tagged with a unique timestamp (version) to avoid ABA problem.

struct

| `value: uint32_t`

| `version: uint32_t`

end

`timestamp_t` = The type of the timestamp tagged with a unique timestamp (version) to avoid ABA problem.

struct

| `value: uint32_t`

```
| | version: uint32_t
| end
node_t = The type of a tree node.
| struct
| | rank: rank_t
| end
```

The shared variables in our LTQueue version are as follows.

Note that we have described a very specific and simple way to organize the tree nodes in dLTQueue in a min-heap-like array structure hosted on the sole dequeuer. We will resume our description of the related tree-structure procedures `parent()` (Procedure 8), `children()` (Procedure 9), `leafNodeIndex()` (Procedure 10) with this representation in mind. However, our algorithm does not strictly require this representation and can be substituted with other more-optimized representations & distributed placements, as long as the similar tree-structure procedures are supported.

Shared variables

```
Counter: gptr<uint64_t>
| A distributed counter shared by the enqueueers. Hosted at the dequeuer.
Tree_size: uint64_t
| A read-only variable storing the number of tree nodes present in the dLTQueue.
Nodes: gptr<node_t>
| An array with Tree_size entries storing all the tree nodes present in the
| dLTQueue shared by all processes.
| Hosted at the dequeuer.
| This array is organized in a similar manner as a min-heap: At index 0 is the root
| node. For every index  $i > 0$ ,  $\lfloor \frac{i-1}{2} \rfloor$  is the index of the parent of node  $i$ . For every
| index  $i > 0$ ,  $2i + 1$  and  $2i + 2$  are the indices of the children of node  $i$ .
Dequeuer_rank: uint32_t
| The rank of the dequeuer process. This is read-only.
Timestamps: A read-only array  $[0..size - 1]$  of gptr<timestamp_t>, with size
| being the number of processes.
| The entry at index  $i$  corresponds to the Min_timestamp distributed variable at the
| process with a rank of  $i$ .
```

Enqueueer-local variables

```
Process_count: uint64_t
| The number of processes.
Self_rank: uint32_t
```

```
| The rank of the current enqueueer
| process.
Min_timestamp: gptr<timestamp_t>
Spssc: spsc_t
```

This SPSC is synchronized with the dequeuer.

Dequeuer-local variables

Process_count: uint64_t

| The number of processes.

Spsc: An **array** of spsc_t with Process_count entries.

| The entry at index i corresponds to the Spsc at the process with a rank of i .

Initially, the enqueueers and the dequeuer are initialized as follows:

Enqueuer initialization

Initialize Process_count, Self_rank and Dequeuer_rank.

Initialize Spsc to the initial state.

Initialize Min_timestamp to timestamp_t {MAX_TIMESTAMP, 0}.

Dequeuer initialization

Initialize Process_count, Self_rank and Dequeuer_rank.

Initialize Counter to 0.

Initialize Tree_size to Process_count * 2.

Initialize Nodes to an array with Tree_size entries. Each entry is initialized to node_t {DUMMY_RANK}.

Initialize Spsc, synchronizing each entry with the corresponding enqueuer.

Initialize Timestamps, synchronizing each entry with the corresponding enqueuer.

4.2.3 Algorithm

We first present the tree-structure utility procedures that are shared by both the enqueuer and the dequeuer:

Procedure 8: uint32_t parent(uint32_t index)

```
1 return (index - 1) / 2
```

parent returns the index of the parent tree node given the node with index index. These indices are based on the shared Nodes array. Based on how we organize the Nodes array, the index of the parent tree node of index is $(\text{index} - 1) / 2$.

Procedure 9: `vector<uint32_t> children(uint32_t index)`

```
2 left_child = index * 2 + 1
3 right_child = left_child + 1
4 res = vector<uint32_t>()
5 if (left_child >= Tree_size)
6 | return res
7 res.push(left_child)
8 if (right_child >= Tree_size)
9 | return res
10 res.push(right_child)
11 return res
```

Similarly, `children` returns all indices of the child tree nodes given the node with index `index`. These indices are based on the shared `Nodes` array. Based on how we organize the `Nodes` array, these indices can be either `index * 2 + 1` or `index * 2 + 2`.

Procedure 10: `uint32_t leafNodeIndex(uint32_t enqueue_rank)`

```
12 return Tree_size + enqueue_rank
```

`leafNodeIndex` returns the index of the leaf node that is logically attached to the enqueue node with rank `enqueue_rank` as in Figure 19.

The followings are the enqueue procedures.

Procedure 11: `bool enqueue(data_t value)`

```
13 timestamp = faa(Counter, 1)
14 if (!spsc_enqueue(&Spsc, (value, timestamp)))
15 | return false
16 propagate_e()
17 return true
```

To enqueue a value, `enqueue` first obtains a count by FAA-ing the distributed counter `Counter` (Line 13). Then, we enqueue the data tagged with the timestamp into the local SPSC (Line 14). Then, enqueue propagates the changes by invoking `propagate_e()` (Line 16) and returns `true`.

Procedure 12: void propagate_e()

```
18 if (!refreshTimestampe())
19   | refreshTimestampe()
20 if (!refreshLeafe())
21   | refreshLeafe()
22 current_node_index = leafNodeIndex(Self_rank)
23 repeat
24   | current_node_index = parent(current_node_index)
25   | if (!refreshe(current_node_index))
26     | refreshe(current_node_index)
27 until current_node_index == 0
```

The propagate_e procedure is responsible for propagating SPSC updates up to the root node as a way to notify other processes of the newly enqueued item. It is split into 3 phases: Refreshing of Min_timestamp in the enqueuer node (Line 18 - Line 19), refreshing of the enqueuer's leaf node (Line 20 - Line 21), refreshing of internal nodes (Line 23 - Line 27). On Line 20 - Line 27, we refresh every tree node that lies between the enqueuer node and the root node.

Procedure 13: bool refreshTimestamp_e()

```
28 min_timestamp = timestamp_t {}
29 read(Min_timestamp, &min_timestamp)
30 {old_timestamp, old_version} = min_timestamp
31 front = (data_t {}, timestamp_t {})
32 is_empty = !spsc_readFront(Spsc, &front)
33 if (is_empty)
34   | return cas(Min_timestamp,
35               timestamp_t {old_timestamp, old_version},
36               timestamp_t {MAX_TIMESTAMP, old_version + 1})
35 else
36   | return cas(Min_timestamp,
37               timestamp_t {old_timestamp, old_version},
38               timestamp_t {front.timestamp, old_version + 1})
```

The refreshTimestamp_e procedure is responsible for updating the Min_timestamp of the enqueuer node. It simply looks at the front of the local SPSC (Line 32) and CAS Min_timestamp accordingly (Line 33 - Line 36).

Procedure 14: `bool refreshNodee(uint32_t current_node_index)`

```
37 current_node = node_t {}
38 read(Nodes, current_node_index, &current_node)
39 {old_rank, old_version} = current_node.rank
40 min_rank = DUMMY_RANK
41 min_timestamp = MAX_TIMESTAMP
42 for child_node_index in children(current_node)
43     child_node = node_t {}
44     read(Nodes + child_node_index, &child_node)
45     {child_rank, child_version} = child_node
46     if (child_rank == DUMMY_RANK) continue
47     child_timestamp = timestamp_t {}
48     read(Timestamps + child_rank, &child_timestamp)
49     if (child_timestamp < min_timestamp)
50         min_timestamp = child_timestamp
51         min_rank = child_rank
    return cas(Nodes + current_node_index,
52 node_t {rank_t {old_rank, old_version}},
    node_t {rank_t {min_rank, old_version + 1}})
```

The `refreshNodee` procedure is responsible for updating the ranks of the internal nodes affected by the enqueue. It loops over the children of the current internal nodes (Line 42). For each child node, we read the rank stored in it (Line 45), if the rank is not `DUMMY_RANK`, we proceed to read the value of `Min_timestamp` of the enqueuer node with the corresponding rank (Line 48). At the end of the loop, we obtain the rank stored inside one of the child nodes that has the minimum timestamp stored in its enqueuer node (Line 50 - Line 51). We then try to CAS the rank inside the current internal node to this rank (Line 52).

Procedure 15: `bool refreshLeafe()`

```
53 leaf_node_index = leafNodeIndex(Self_rank)
54 leaf_node = node_t {}
55 read(Nodes + leaf_node_index, &leaf_node)
56 {old_rank, old_version} = leaf_node.rank
57 min_timestamp = timestamp_t {}
58 read(Min_timestamp, &min_timestamp)
59 timestamp = min_timestamp.timestamp
    return cas(Nodes + leaf_node_index,
60 node_t {rank_t {old_rank, old_version}},
    node_t {timestamp == MAX ? DUMMY_RANK : Self_rank, old_version + 1})
```

The `refreshLeafe` procedure is responsible for updating the rank of the leaf node affected by the enqueue. It simply reads the value of `Min_timestamp` of the enqueueer node it is logically attached to (Line 58) and CAS the leaf node's rank accordingly (Line 60).

The followings are the dequeuer procedures.

Procedure 16: `bool dequeue(data_t* output)`

```
61 root_node = node_t {}
62 read(Nodes, &root_node)
63 {rank, version} = root_node.rank
64 if (rank == DUMMY_RANK) return false
65 output_with_timestamp = (data_t {}, timestamp_t {})
66   if (!spsc_dequeue(&Spscs[rank]),
67       &output_with_timestamp))
67   | return false
68 *output = output_with_timestamp.data
69 propagated(rank)
70 return true
```

To dequeue a value, `dequeue` reads the rank stored inside the root node (Line 63). If the rank is `DUMMY_RANK`, the MPSC queue is treated as empty and failure is signaled (Line 64). Otherwise, we invoke `spsc_dequeue` on the SPSC of the enqueueer with the obtained rank (Line 66). We then extract out the real data and set it to output (Line 68). We finally propagate the dequeue from the enqueueer node that corresponds to the obtained rank (Line 69) and signal success (Line 70).

Procedure 17: `void propagated(uint32_t enqueueer_rank)`

```
71 if (!refreshTimestampd(enqueueer_rank))
72   | refreshTimestampd(enqueueer_rank)
73 if (!refreshLeafd(enqueueer_rank))
74   | refreshLeafd(enqueueer_rank)
75 current_node_index = leafNodeIndex(enqueueer_rank)
76 repeat
77   | current_node_index = parent(current_node_index)
78   | if (!refreshd(current_node_index))
79     | refreshd(current_node_index)
80 until current_node_index == 0
```

The `propagated` procedure is similar to `propagatee`, with appropriate changes to accommodate the dequeuer.

Procedure 18: `bool refreshTimestampd(uint32_t enqueue_rank)`

```
81 min_timestamp = timestamp_t {}
82 read(Timestamps + enqueue_rank, &min_timestamp)
83 {old_timestamp, old_version} = min_timestamp
84 front = (data_t {}, timestamp_t {})
85 is_empty = !spsc_readFront(&Spscs[enqueue_rank], &front)
86 if (is_empty)
87     return cas(Timestamps + enqueue_rank,
88               timestamp_t {old_timestamp, old_version},
89               timestamp_t {MAX_TIMESTAMP, old_version + 1})
88 else
89     return cas(Timestamps + enqueue_rank,
90               timestamp_t {old_timestamp, old_version},
91               timestamp_t {front.timestamp, old_version + 1})
```

The `refreshTimestampd` procedure is similar to `refreshTimestampe`, with appropriate changes to accommodate the dequeuer.

Procedure 19: `bool refreshNoded(uint32_t current_node_index)`

```
91 current_node = node_t {}
92 read(Nodes + current_node_index, &current_node)
93 {old_rank, old_version} = current_node.rank
94 min_rank = DUMMY_RANK
95 min_timestamp = MAX_TIMESTAMP
96 for child_node_index in children(current_node)
97     child_node = node_t {}
98     read(Nodes + child_node_index, &child_node)
99     {child_rank, child_version} = child_node
100     if (child_rank == DUMMY_RANK) continue
101     child_timestamp = timestamp_t {}
102     read(Timestamps + child_rank, &child_timestamp)
103     if (child_timestamp < min_timestamp)
104         min_timestamp = child_timestamp
105         min_rank = child_rank
106     return cas(Nodes + current_node_index,
107               node_t {rank_t {old_rank, old_version}},
108               node_t {rank_t {min_rank, old_version + 1}})
```

The `refreshNoded` procedure is similar to `refreshNodee`, with appropriate changes to accommodate the dequeuer.

Procedure 20: `bool refreshLeafd(uint32_t enqueue_rank)`

```
107 leaf_node_index = leafNodeIndex(enqueue_rank)
108 leaf_node = node_t {}
109 read(Nodes + leaf_node_index, &leaf_node)
110 {old_rank, old_version} = leaf_node.rank
111 min_timestamp = timestamp_t {}
112 read(Timestamps + enqueue_rank, &min_timestamp)
113 timestamp = min_timestamp.timestamp
    return cas(Nodes + leaf_node_index,
114 node_t {rank_t {old_rank, old_version}},
    node_t {timestamp == MAX ? DUMMY_RANK : Self_rank, old_version + 1})
```

The `refreshLeafd` procedure is similar to `refreshLeafe`, with appropriate changes to accommodate the dequeuer.

4.3 Slotqueue - dLTQueue-inspired distributed MPSC queue with all constant-time operations

The straightforward dLTQueue algorithm we have ported in Section 4.2 pretty much preserves the original algorithm's characteristics, i.e. wait-freedom and time complexity of $\Theta(\log n)$ for dequeue and enqueue operations. We note that in shared-memory systems, this logarithmic growth is fine. However, in distributed systems, this increase in remote operations would present a bottleneck in enqueue and dequeue latency. Upon closer inspection, this logarithmic growth is due to the propagation process because it has to traverse every level in the tree. Intuitively, this is the problem of we trying to maintain the tree structure. Therefore, to be more suitable for distributed context, we propose a new algorithm Slotqueue inspired by LTQueue, which uses a slightly different structure. The key point is that both enqueue and dequeue only perform a constant number of remote operations, at the cost of dequeue having to perform $\Theta(n)$ local operations, where n is the number of enqueueers. Because remote operations are much more expensive, this might be a worthy tradeoff.

4.3.1 Overview

The structure of Slotqueue is shown as in Figure 20.

Each enqueueer hosts a distributed SPSC as in dLTQueue (Section 4.2). The enqueueer when enqueues a value to its local SPSC will timestamp the value using a distributed counter hosted at the dequeuer.

Additionally, the dequeuer hosts an array whose entries each corresponds with an enqueueer. Each entry stores the minimum timestamp of the local SPSC of the corresponding enqueueer.



Figure 20: Basic structure of Slotqueue.

4.3.2 Data structure

We first introduce the types and shared variables utilized in Slotqueue.

Types

`data_t` = The type of data stored.

`timestamp_t` = `uint64_t`

`spsc_t` = The type of the SPSC each enqueueer uses, this is assumed to be the distributed SPSC in Section 4.1.

Shared variables

Slots: `gptr<timestamp_t*>`

An array of `timestamp_t` with the number of entries equal to the number of enqueueers.

Hosted at the dequeueer.

Counter: `gptr<uint64_t>`

A distributed counter.

Hosted at the dequeueer.

Enqueueer-local variables

Dequeueer_rank: `uint64_t`

The rank of the dequeueer.

Process_count: `uint64_t`

| The number of enqueueers.
Self_rank: uint32_t
| The rank of the current enqueueer process.
Spsc: spsc_t
| This SPSC is synchronized with the dequeuer.

Dequeuer-local variables

Dequeuer_rank: uint64_t
| The rank of the dequeuer.
Process_count: uint64_t
| The number of enqueueers.
Spscs: An **array** of spsc_t with Process_count entries.
| The entry at index i corresponds to the Spsc at the process with a rank of i .

Initially, the enqueueer and the dequeuer are initialized as follows.

Enqueueer initialization

Initialize Dequeuer_rank.
Initialize Process_count.
Initialize Self_rank.
Initialize the local Spsc to its initial state.

Dequeuer initialization

Initialize Dequeuer_rank.
Initialize Process_count.
Initialize Counter to 0.
Initialize the Slots array with size equal to the number of enqueueers and every entry is initialized to MAX_TIMESTAMP.
Initialize the Spscs array, the i -th entry corresponds to the Spsc variable of the process of rank i .

4.3.3 Algorithm

The enqueueer operations are given as follows.

Procedure 21: bool enqueue(data_t v)

```
1 timestamp = faa(Counter, 1)
2 if (!spsc_enqueue(&Spsc, (v, timestamp))) return false
3 if (!refreshEnqueue(timestamp))
4 | refreshEnqueue(timestamp)
5 return true
```

To enqueue a value, enqueue first obtains a timestamp by FAA-ing the distributed counter (Line 1). It then tries to enqueue the value tagged with the timestamp (Line 2). At Line 3 - Line 4, the enqueueer tries to refresh its slot's timestamp.

Procedure 22: `bool refreshEnqueue(timestamp_t ts)`

```
6 front = (data_t {}, timestamp_t {})  
7 success = spsc_readFront(Spsc, &front)  
8 new_timestamp = success ? front.timestamp : MAX_TIMESTAMP  
9 if (new_timestamp != ts)  
10 | return true  
11 old_timestamp = timestamp_t {}  
12 read(Slots + Self_rank, &old_timestamp)  
13 success = spsc_readFront(Spsc, &front)  
14 new_timestamp = success ? front.timestamp : MAX_TIMESTAMP  
15 if (new_timestamp != ts)  
16 | return true  
    return cas(Slots + Self_rank,  
17     old_timestamp,  
        new_timestamp)
```

refreshEnqueue's responsibility is to refresh the timestamp stores in the enqueueer's slot to potentially notify the dequeuer of its newly-enqueued element. It first reads the current front element (Line 7). If the SPSC is empty, the new timestamp is set to MAX_TIMESTAMP, otherwise, the front element's timestamp (Line 8). If it finds that the front element's timestamp is different from the timestamp ts it returns true immediately (Line 9 - Line 10). Otherwise, it reads its slot's old timestamp (Line 12) and re-reads the current front element in the SPSC (Line 13) to update the new timestamp. Note that similar to Line 10, refreshEnqueue immediately succeeds if the new timestamp is different from the timestamp ts of the element it enqueues (Line 16). Otherwise, it tries to CAS its slot's timestamp with the new timestamp (Line 17).

The dequeuer operations are given as follows.

Procedure 23: `bool dequeue(data_t* output)`

```
19 rank = readMinimumRank()  
20 if (rank == DUMMY_RANK)  
21 | return false  
22 output_with_timestamp = (data_t {}, timestamp_t {})  
23 if (!spsc_dequeue(Spsc, &output_with_timestamp))  
24 | return false  
25 *output = output_with_timestamp.data  
26 if (!refreshDequeue(rank))  
27 | refreshDequeue(rank)  
28 return true
```

To dequeue a value, dequeue first reads the rank of the enqueueer whose slot currently stores the minimum timestamp (Line 19). If the obtained rank is DUMMY_RANK, failure is signaled (Line 20 - Line 21). Otherwise, it tries to dequeue the SPSC of the corresponding enqueueer (Line 23). It then tries to refresh the enqueueer's slot's timestamp to potentially notify the enqueueer of the dequeue (Line 26 - Line 27). It then signals success (Line 28).

Procedure 24: uint64_t readMinimumRank()

```
29 buffered_slots = timestamp_t[Process_count] {}
30 for index in 0..Process_count
31   | read(Slots + index, &buffered_slots[index])
32 if every entry in buffered_slots is MAX_TIMESTAMP
33   | return DUMMY_RANK
34   let rank be the index of the first slot that contains the minimum timestamp
      among buffered_slots
35 for index in 0..rank
36   | read(Slots + index, &buffered_slots[index])
37 min_timestamp = MAX_TIMESTAMP
38 for index in 0..rank
39   | timestamp = buffered_slots[index]
40   | if (min_timestamp < timestamp)
41     | min_rank = index
42     | min_timestamp = timestamp
43 return min_rank
```

readMinimumRank's main responsibility is to return the rank of the enqueueer from which we can safely dequeue next. It first creates a local buffer to store the value read from Slots (Line 29). It then performs 2 scans of Slots and read every entry into buffered_slots (Line 30 - Line 36). If the first scan finds only MAX_TIMESTAMPS, DUMMY_RANK is returned (Line 33). From there, based on buffered_slots, it returns the rank of the enqueueer whose buffered slot stores the minimum timestamp (Line 38 - Line 43).

Procedure 25: refreshDequeue(rank: int) **returns** bool

```
46 old_timestamp = timestamp_t {}
47 read(Slots + rank, &old_timestamp)
48 front = (data_t {}, timestamp_t {})
49 success = spsc_readFront(&Spscs[rank], &front)
50 new_timestamp = success ? front.timestamp : MAX_TIMESTAMP
   return cas(Slots + rank,
51     old_timestamp,
       new_timestamp)
```

refreshDequeue's responsibility is to refresh the timestamp of the just-dequeued enqueueer to notify the enqueueer of the dequeue. It first reads the old timestamp of the slot (Line 47) and the front element (Line 49). If the SPSC is empty, the new timestamp is set to MAX_TIMESTAMP, otherwise, it is the front element's timestamp (Line 50). It finally tries to CAS the slot with the new timestamp (Line 51).

Chapter V Evaluation

This section introduces our benchmarking process, including our setup, environment, metrics of interest, and our microbenchmark program. Most importantly, we showcase the preliminary results on how well our novel algorithms perform, especially Slotqueue. We conclude this section with a discussion about the implications of these results.

Currently, performance-related properties are our main focus.

5.1 Benchmarking metrics

This section provides an overview of the metrics we are interested in for our algorithms. Performance-wise, latency and throughput are the two most popular metrics. These metrics revolve around the concept of a “task”. In our context, a task is a single method call of an MPSC queue algorithm, e.g., enqueue and dequeue. Note that in our discussion, any two tasks are independent. Roughly speaking, two tasks are independent if one does not need to depend on the output of another for it to finish or there does not exist a bigger task that needs to depend on the outputs of the tasks. This rules out pipeline parallelism, where a task needs to wait for the output of a preceding task, and data parallelism, where a big task is split into and needs to wait for the outputs of multiple smaller tasks.

5.1.1 Throughput

Throughput is the number of operations finished in a unit of time. Its unit is often given as $\frac{\text{ops}}{\text{s}}$ (operations per second), $\frac{\text{ops}}{\text{ms}}$ (operations per millisecond), or $\frac{\text{ops}}{\text{us}}$ (operations per microsecond). Intuitively, throughput is closest to our notion of “performance”: The higher the throughput, the more tasks are done in a unit of time and, thus, the higher the performance. The implication is that our ultimate goal is to optimize the throughput metric of our algorithms.

5.1.2 Latency

Latency is the time it takes for a single task to complete. Its unit is often given as $\frac{\text{s}}{\text{op}}$ (seconds per operation), $\frac{\text{ms}}{\text{op}}$ (milliseconds per operation), or $\frac{\text{us}}{\text{op}}$ (microseconds per operation).

Intuitively, to optimize latency, one should minimize the number of execution steps required by a task. Therefore, it is obvious that optimizing for latency is much clearer than optimizing for throughput.

In concurrent algorithms, multiple tasks are executed by multiple processes. We observe that, if we fix the number of processes, the lower the average latency of a task, the larger the number of tasks that can be completed by a process, which implies higher throughput. Therefore, good latency often (but not always) implies good throughput.

From the two points above, we can see that latency is a more intuitive metric to optimize for, while being quite indicative of the algorithm's performance.

One question is: how do we optimize for latency? As we have discussed, we should minimize the number of execution steps. A key observation is that when the number of processes grows, contention should also grow, thus causing the number of steps taken by a task to grow and, thus, the average latency to deteriorate. Note that if we manage to keep the average latency of a task fixed while also increasing the number of processes, we gain higher throughput due to higher concurrency. The actionable insight is that if we minimize contention in our algorithms, our algorithms should scale with the number of processes.

Following this discussion, we should aim to discover and optimize highly contended areas in our algorithms if we want to make them scale well to a large number of nodes/processes.

5.2 Benchmarking baselines

We use three MPSC queue algorithms as benchmarking baselines:

- dLTQueue + our custom SPSC: Our most optimized version of LTQueue while still keeping the core algorithm intact.
- Slotqueue + our custom SPSC: Our modification to dLTQueue to obtain a more optimized distributed version of LTQueue.
- AMQueue [1]: A hosted bounded MPSC queue algorithm, already detailed in Section 3.2.

5.3 Microbenchmark program

Our microbenchmark is as follows:

- All processes share a single MPSC; one of the processes is a dequeuer, and the rest are enqueueers.
- The enqueueers enqueue a total of 10^4 elements.
- The dequeuer dequeues 10^4 elements.
- For MPSC, the MPSC is warmed up before the dequeuer starts.

We measure the latency and throughput of the enqueue and dequeue operations. This microbenchmark is repeated 5 times for each algorithm, and we take the mean of the results.

5.4 Benchmarking setup

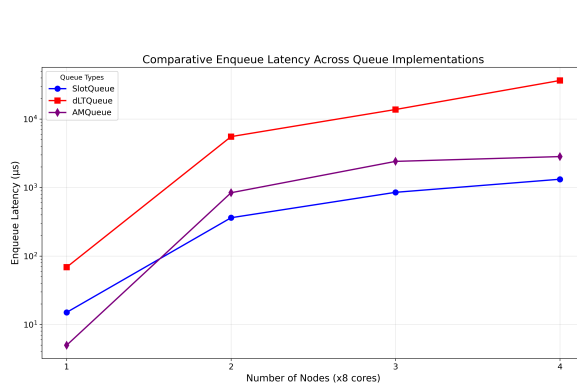
The experiments are carried out on a four-node cluster residing in the HPC Lab at Ho Chi Minh University of Technology. Each node is an Intel Xeon CPU E5-2680 v3, which has 8 cores and 16 GB RAM. The interconnect used is Ethernet and, thus, does not support true one-sided communication.

The operating system used is Ubuntu 22.04.5. The MPI implementation used is MPICH version 4.0, released on January 21, 2022.

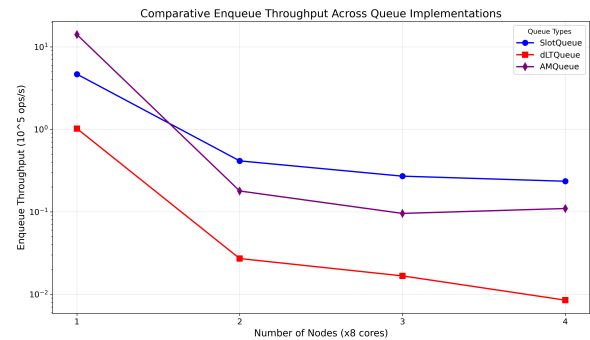
We run the producer-consumer microbenchmark on 1 to 4 nodes to measure both the latency and performance of our MPSC algorithms.

5.5 Benchmarking results

Figure 21, Figure 22, and Figure 23 showcase our benchmarking results, with the y-axis drawn in log scale.

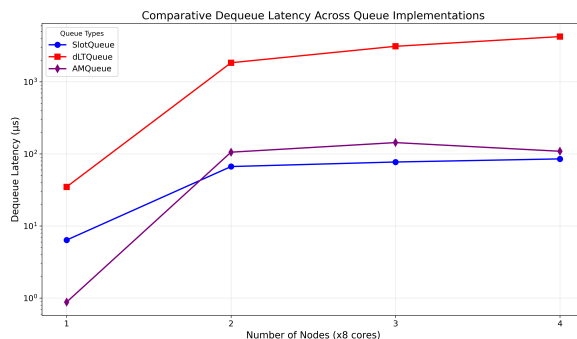


(a) Enqueue latency benchmark results.

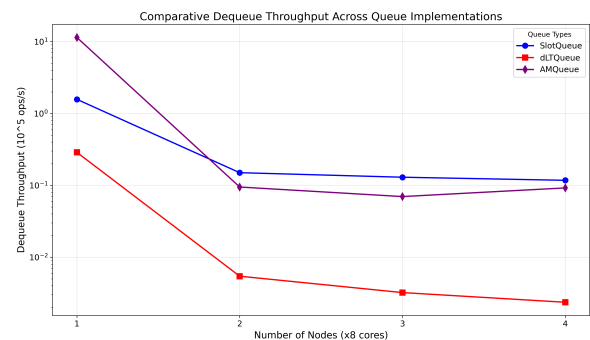


(b) Enqueue throughput benchmark results

Figure 21: Microbenchmark results for enqueue operation.



(a) Dequeue latency benchmark results.



(b) Dequeue throughput benchmark results

Figure 22: Microbenchmark results for dequeue operation.

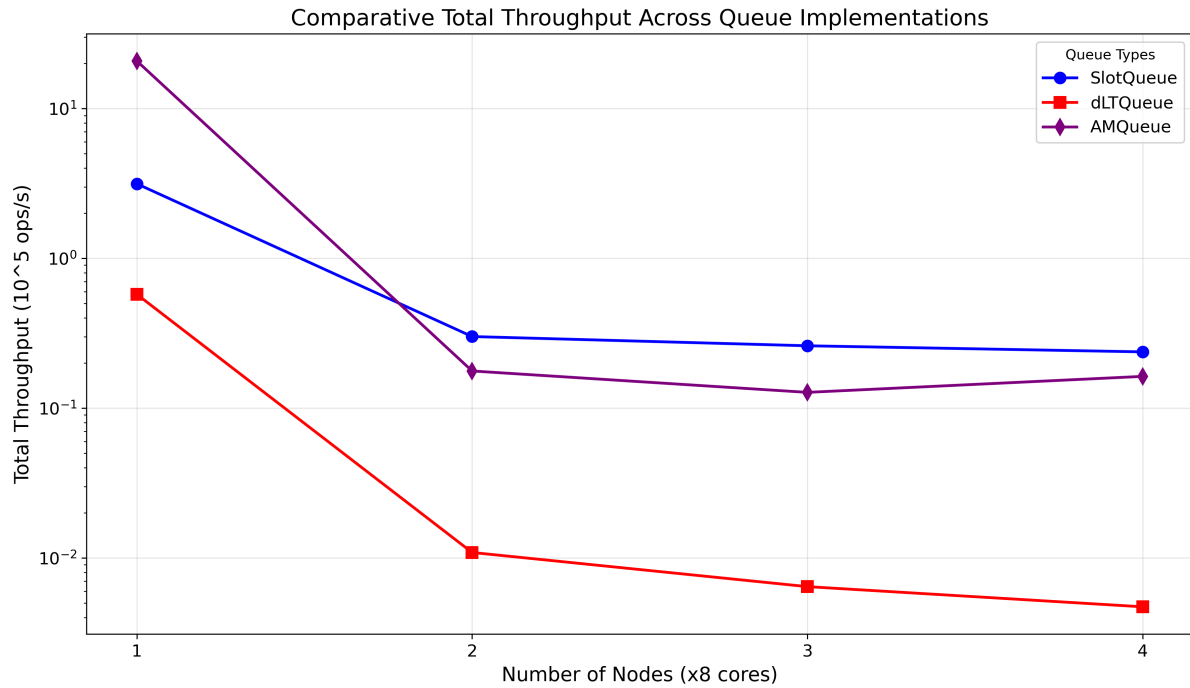


Figure 23: Microbenchmark results for total throughput.

The most evident thing is that Figure 23 and Figure 22b are almost identical. This supports our claim that in an MPSC queue, the performance is bottlenecked by the dequeuer.

For enqueue latency and throughput, Slotqueue performs far better than dLTQueue while being slightly better than AMQueue. This is in line with our theoretical projection in Table 4. One concerning trend is that Slotqueue's enqueue throughput seems to degrade with the number of nodes, which signals a potential scalability problem. This is further problematic in that our theoretical model suggests that the cost of enqueue is always fixed. This is to be investigated further in the future.

For dequeue latency and throughput, Slotqueue and AMQueue are quite closely matched, while being better than dLTQueue. This is expected, agreeing with our projection of dequeue wrapping overhead in Table 4. Furthermore, Slotqueue is conceived as a more dequeuer-optimized version of dLTQueue. Based on this empirical result, it is reasonable to believe this to be the case. Unlike enqueue, the dequeue latency of Slotqueue seems to be quite stable, increasing very slowly. Because the dequeuer is the bottleneck of an MPSC, this is a good sign for the scalability of Slotqueue.

In conclusion, based on Figure 23, Slotqueue seems to perform better than dLTQueue and AMQueue in terms of both enqueue and dequeue operations, both latency-wise and throughput-wise. The overhead of a logarithmic-order number of remote operations in dLTQueue seems to be costly, adversely affecting its performance when the number of nodes increases. Additionally, compared to AMQueue, dLTQueue and Slotqueue also have the advantage of fault tolerance, which, due to the blocking nature of AMQueue, cannot be promised.

Chapter VI Conclusion

This thesis explores how shared-memory programming principles, particularly the use of atomic operations, can be applied to model and develop distributed MPSC queue algorithms. We examined existing MPSC queue algorithms from shared-memory literature and adapted them for distributed computing environments using our proposed model. Based on this foundation, we developed two novel distributed MPSC queue algorithms: dLTQueue and Slotqueue. We provided formal proofs demonstrating key theoretical properties of these algorithms, including their correctness, fault-tolerance capabilities, and performance characteristics. To validate our theoretical findings, we conducted empirical benchmarks comparing queue behavior with an existing algorithm called Active Message Queue (AMQueue) referenced in [1].

Appendix A Theoretical aspects

This section discusses the correctness and progress guarantee properties of the distributed MPSC queue algorithms introduced in Chapter IV. We also provide a theoretical performance model of these algorithms to predict how well they scale to multiple nodes.

A.1 Preliminaries

In this section, we formalize the notion of correct concurrent algorithms (Section A.1.1). Verifying the correctness of concurrent queues involves us giving a sequential queue specification based on the Larch interface specification [25] (Section A.1.2), Herlihy and Wing's method for verifying linearizability [26] (Section A.1.3), the Owicki-Gries formal method for reasoning about concurrent algorithms [27] (Section A.1.4). Finally, we formulate the definition of harmless ABA problem in Section A.1.5.

We will base our proofs on these formalisms to prove the algorithms' correctness.

Our system consists of a set of sequential processes that communicate through a collection of shared objects. Processes are asynchronous, so that each process may run at their own pace. Each object has a type, defining a set of possible values and operations that manipulate that object.

A.1.1 Linearizability

This section provides the formal definition of linearizability, which was not given in Section 2.3. Our formalism is based on Herlihy and Wing's notion introduced in [21] and [26]. Specification and verification of linearizable objects (queues) will be deferred to Section A.1.2 and Section A.1.3.

An execution of a concurrent system is modeled by a history, which is a finite sequence of operation *invocation* and *response events* [26]:

- An invocation is of the form $x \text{ op}(\text{args}^*) \ A$ where x is the object names, op is the operation name, arg^* is the list of arguments and A is the name of a process.
- A response is of the form $x \text{ term}(\text{res}^*) \ A$ where x is the object names, term is the termination status (which is assumed to be Ok in this thesis for normal termination), res^* is the list of results and A is the name of a process.

A response event *matches* an invocation event if their object names and process names are the same. If there is no matching response event for an invocation event, the response event is said to be *pending*. $\text{Complete}(H)$ is a history obtained from a history H by removing all pending events in it.

A history is *sequential* when it begins with an invocation event, and every invocation event is paired with a corresponding response event that follows it (with the exception that the final invocation may not yet have its response).

Given a history H , we can extract two types of subsequences:

- A *process subhistory* $H|P$ contains only those events from H where the process name matches P .
- An *object subhistory* $H \upharpoonright x$ contains only those events from H where the object name matches x .

Two histories H and H' are considered equivalent when their process subhistories are identical for every process P .

We assume all history to be *well-formed*, that is the history H such that $H|P$ is sequential for every P .

An *operation* e within a history is defined as a pair composed of an invocation $inv(e)$ and the subsequent matching response $res(e)$. Operation e_0 *lies within* operation e_1 in history H if e_1 's invocation comes first, then e_0 's invocation, then e_0 's response, and finally e_1 's response. Operation e_0 *precedes* operation e_1 if e_0 's response comes before e_1 's invocation. A history H induces a precedence strict partial order \prec_H on operations. That is, $e_0 \prec_H e_1$ iff e_0 precedes e_1 .

Definition A.1.1.1 (*Linearizability* [26]) A history H is *linearizable* if can be extended (by appending zero or more events) to some history H' such that:

- $Complete(H')$ is equivalent to some legal sequential history S .
- $\prec_{H'} \subseteq \prec_S$.

In this case S is called a *linearization* of H .

A.1.2 Larch interface specification of queues

A.1.3 Verifying linearizability

A.1.4 The Owicki-Gries method

A.1.5 ABA-safety

A.2 Theoretical proofs of the distributed SPSC

This section establishes the correctness and fault-tolerance of our simple distributed SPSC introduced in Section 4.1. Specifically, Section A.2.1 proves there's no ABA problem and unsafe memory reclamation in our SPSC queue and it is linearizable; Section A.2.2 proves that our SPSC queue is wait-free; Section A.2.3 discusses the overhead involved in each SPSC operation.

A.2.1 Correctness

First, we prove that there's no ABA problem in our SPSC queue in Section A.2.1.1. Second, we prove that there's no potential memory errors when executing our SPSC

queue in Section A.2.1.2. Finally and most importantly, we prove that our SPSC queue is linearizable in Section A.2.1.3.

A.2.1.1 ABA problem

There's no CAS instruction in our simple distributed SPSC, so there's no potential for ABA problem.

A.2.1.2 Memory reclamation

There's no dynamic memory allocation and deallocation in our simple distributed SPSC, so it is memory-safe.

A.2.1.3 Linearizability

A.2.2 Progress guarantee

Our simple distributed SPSC is wait-free:

- `spsc_dequeue` (Procedure 6) does not execute any loops or wait for any other method calls.
- `spsc_enqueue` (Procedure 4) does not execute any loops or wait for any other method calls.
- `spsc_readFronte` (Procedure 5) does not execute any loops or wait for any other method calls.
- `spsc_readFrontd` (Procedure 7) does not execute any loops or wait for any other method calls.

A.2.3 Theoretical performance

A summary of the theoretical performance of our simple SPSC is provided in Table 5. In the following discussion, R means remote operations and L means local operations.

Operations	Time-complexity
<code>spsc_enqueue</code>	$R + L$
<code>spsc_dequeue</code>	$R + L$
<code>spsc_readFront_e</code>	$R + L$
<code>spsc_readFront_d</code>	R

Table 5: Theoretical performance summary of our simple distributed SPSC. R means remote operations and L means local operations.

For `spsc_enqueue`, we consider the procedure Procedure 4. In the usual case, the remote operation on Line 3 is skipped and so only 2 remote puts are performed on Line 6 and Line 7. The Data array on Line 6 is hosted on the enqueueer, so this is actually a local operation, while the control variable is hosted on the dequeuer, so Line 7 is truly a remote operation. Therefore, theoretically, it's one remote operation plus a local one.

For `spsc_dequeue`, we consider the procedure Procedure 6. Similarly, in the usual case, the remote operation on Line 19 is skipped and only the 2 lines Line 22 and Line 23 are executed always. Here, it's the other way around, the access to the Data array on Line 22 is a truly remote operation while the access to the First control variable is a local one. Therefore, theoretically, it's one remote operation plus a local one.

For `spsc_readFronte`, we consider the procedure Procedure 5. The operation on Line 12 is a truly remote operation, as the First control variable is hosted on the dequeuer. The operation on Line 15 is a remote operation, as the Data array is hosted on the enqueuer. This means, theoretically, it also takes one remote operation plus a local one.

For `spsc_readFrontd`, we consider the procedure Procedure 7. Only the operation on Line 30 is executed always, which results in a truly remote operation as the Data array is hosted on the enqueuer. Therefore, it only takes one remote operation.

A.3 Theoretical proofs of dLTQueue

A.3.1 Proof-specific notations

The structure of dLTQueue is presented again in Figure 24.

As a reminder, the bottom rectangular nodes are called the **enqueuer nodes** and the circular node are called the **tree nodes**. Tree nodes that are attached to an enqueuer node are called **leaf nodes**, otherwise, they are called **internal nodes**. Each **enqueuer node** is hosted on the enqueuer that corresponds to it. The enqueuer nodes accommodate an instance of our distributed SPSC in Section 4.1 and a `Min_timestamp` variable representing the minimum timestamp inside the SPSC. Each **tree node** stores a rank of an enqueuer that's attached to the subtree which roots at the **tree node**.

We will refer `propagatee` and `propagated` as `propagate` if there's no need for discrimination. Similarly, we will sometimes refer to `refreshNodee` and `refreshNoded` as `refreshNode`, `refreshLeafe` and `refreshLeafd` as `refreshLeaf`, `refreshTimestampe` and `refreshTimestampd` as `refreshTimestamp`.

Definition A.3.1.1 For a tree node n , the rank stored in n at time t is denoted as $rank(n, t)$.

Definition A.3.1.2 For an enqueue or a dequeue op, the rank of the enqueuer it affects is denoted as $rank(op)$.

Definition A.3.1.3 For an enqueuer whose rank is r , the `Min_timestamp` value stored in its enqueuer node at time t is denoted as $min-ts(r, t)$. If r is `DUMMY_RANK`, $min-ts(r, t)$ is `MAX_TIMESTAMP`.

Definition A.3.1.4 For an enqueuer with rank r , the minimum timestamp among the elements between First and Last in its SPSC at time t is denoted as $min-spsc-ts(r, t)$. If r is dummy, $min-spsc-ts(r, t)$ is `MAX`.

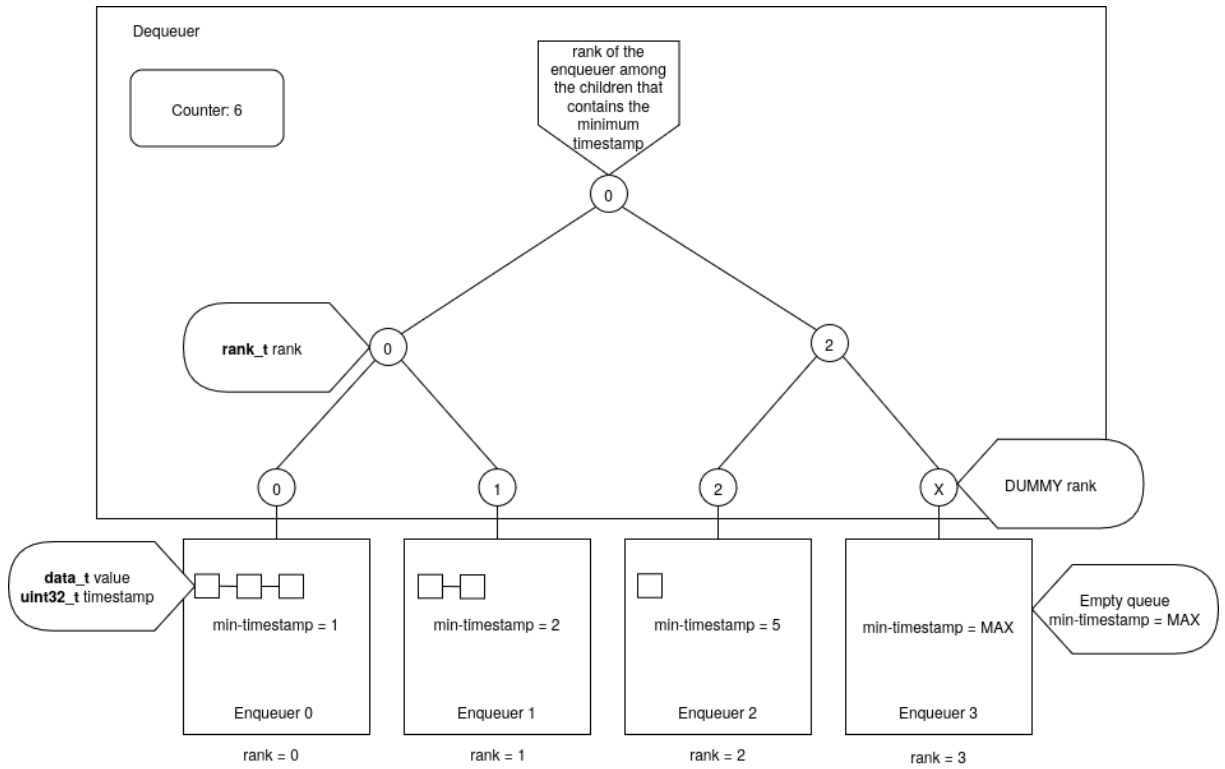


Figure 24: dLTQueue's structure.

Definition A.3.1.5 For an enqueue or a dequeue op, the set of nodes that it calls refreshNode (Procedure 14 or Procedure 19) or refreshLeaf (Procedure 15 or Procedure 20) on is denoted as $path(op)$.

Definition A.3.1.6 For an enqueue or a dequeue, **timestamp-refresh phase** refer to its execution of Line 18 - Line 19 in $propagate_e$ (Procedure 12) or Line 71 - Line 72 in $propagate_d$ (Procedure 17).

Definition A.3.1.7 For an enqueue op, and a node $n \in path(op)$, **node- n -refresh phase** refer to its execution of:

- Line 20 - Line 21 of $propagate_e$ (Procedure 12) if n is a leaf node.
- Line 25 - Line 26 of $propagate_e$ (Procedure 12) to refresh n 's rank if n is a non-leaf node.

Definition A.3.1.8 For a dequeue op, and a node $n \in path(op)$, **node- n -refresh phase** refer to its execution of:

- Line 73 - Line 74 of $propagate_d$ (Procedure 17) if n is a leaf node.
- Line 78 - Line 79 of $propagate_d$ (Procedure 17) to refresh n 's rank if n is a non-leaf node.

Definition A.3.1.9 $refreshTimestamp_e$ (Procedure 13) is said to start its **CAS-sequence** if it finishes Line 29. $refreshTimestamp_e$ is said to end its **CAS-sequence** if it finishes Line 34 or Line 36.

Definition A.3.1.10 $\text{refreshTimestamp}_d$ (Procedure 18) is said to start its **CAS-sequence** if it finishes Line 82. $\text{refreshTimestamp}_d$ is said to end its **CAS-sequence** if it finishes Line 87 or Line 89.

Definition A.3.1.11 refreshNode_e (Procedure 14) is said to start its **CAS-sequence** if it finishes Line 38. refreshNode_e is said to end its **CAS-sequence** if it finishes Line 52.

Definition A.3.1.12 refreshNode_d (Procedure 19) is said to start its **CAS-sequence** if it finishes Line 92. refreshNode_d is said to end its **CAS-sequence** if it finishes Line 106.

Definition A.3.1.13 refreshLeaf_e (Procedure 15) is said to start its **CAS-sequence** if it finishes Line 55. refreshLeaf_e is said to end its **CAS-sequence** if it finishes Line 60.

Definition A.3.1.14 refreshLeaf_d (Procedure 20) is said to start its **CAS-sequence** if it finishes Line 109. refreshLeaf_d is said to end its **CAS-sequence** if it finishes Line 114.

A.3.2 Correctness

A.3.2.1 ABA problem

We use CAS instructions on:

- Line 34 and Line 36 of $\text{refreshTimestamp}_e$ (Procedure 13).
- Line 52 of refreshNode_e (Procedure 14).
- Line 60 of refreshLeaf_e (Procedure 15).
- Line 87 and Line 89 of $\text{refreshTimestamp}_d$ (Procedure 18).
- Line 106 of refreshNode_d (Procedure 19).
- Line 114 of refreshLeaf_d (Procedure 20).

Notice that at these locations, we increase the associated version tags of the CAS-ed values. These version tags are 32-bit in size, therefore, practically, ABA problem can't virtually occur. It's safe to assume that there's no ABA problem in dLTQueue .

A.3.2.2 Memory reclamation

Notice that dLTQueue pushes the memory reclamation problem to the underlying SPSC. If the underlying SPSC is memory-safe, dLTQueue is also memory-safe.

A.3.2.3 Linearizability

A.3.3 Progress guarantee

Notice that every loop in dLTQueue is bounded, and no method have to wait for another. Therefore, dLTQueue is wait-free.

A.3.4 Theoretical performance

A summary of the theoretical performance of dLTQueue is provided in Table 6, which is already shown in Table 4. In the following discussion, R means remote operations and L means local operations.

Operations	Time-complexity
enqueue	$6 \log_2(n)R + 4 \log_2(n)L$
dequeue	$4 \log_2(n)R + 6 \log_2(n)L$

Table 6: Theoretical performance summary of dLTQueue. R means remote operations and L means local operations.

For enqueue, we consider the procedure Procedure 11. We consider the propagation process, which causes most of the remote operations, while Line 13 and Line 14 are negligible. Notice that the number of node refreshes are proportional to the number of the level of the trees, which is $O(n)$ for n being the number of processes. Each level of the tree in the worst case needs 2 retries, each retry would have to:

- Read the current node (which is a truly remote operation for enqueue).
- Read the two child nodes (which is 2 truly remote operations for enqueue).
- Read the two min-timestamp variables in the two child nodes (which is 2 truly local operations for enqueue).
- Compare-and-swap the current node (which is a truly remote operation for enqueue).

In total, each level requires 6 remote operations and 4 local operations. Therefore, enqueue requires about $6 \log_2(n)R + 4 \log_2(n)L$ operations.

For dequeue, it's similar to enqueue but the other way around, what makes for a remote operation in enqueue is a local operation in dequeue and otherwise. Therefore, dequeue requires about $4 \log_2(n)R + 6 \log_2(n)L$ operations.

A.4 Theoretical proofs of Slotqueue

A.4.1 Proof-specific notations

As a refresher, Figure 25 shows the structure of Slotqueue.

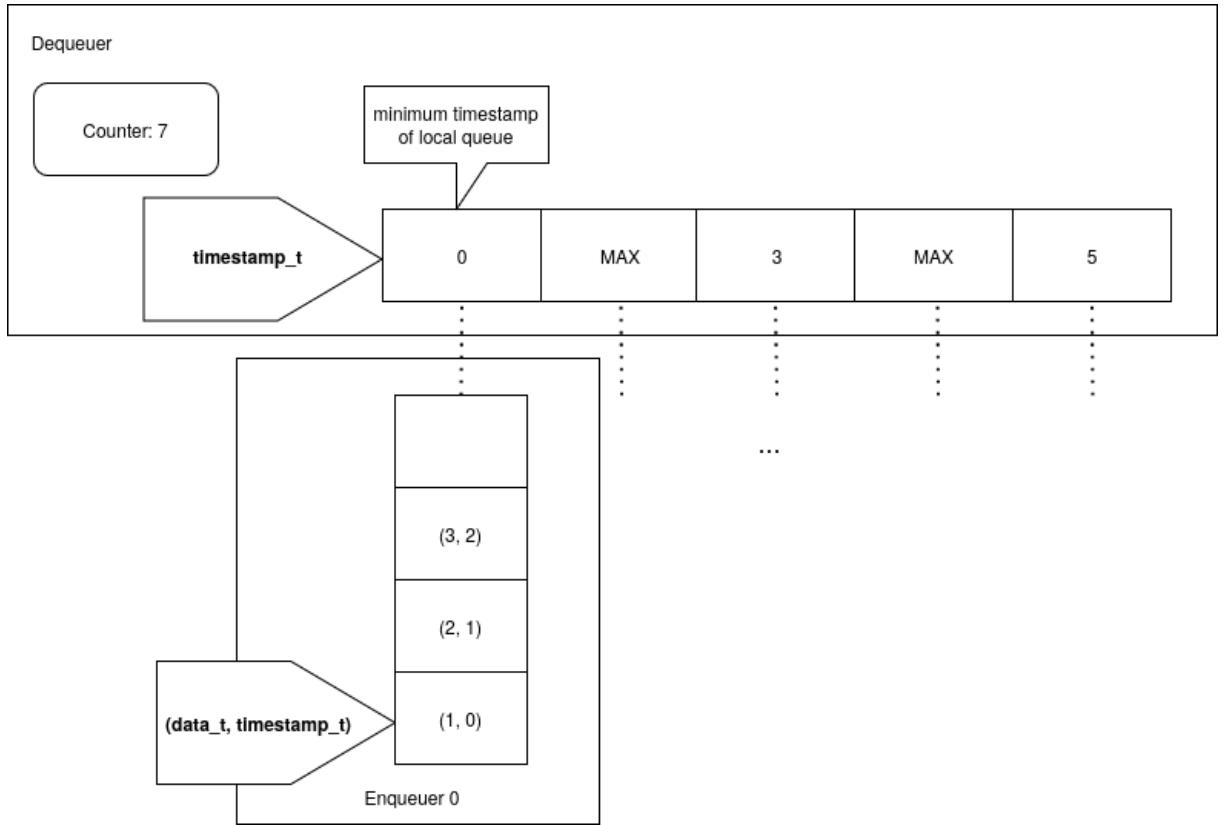


Figure 25: Basic structure of Slotqueue.

Each enqueueer hosts an SPSC that can only be accessed by itself and the dequeueer. The dequeueer hosts an array of slots, each slot corresponds to an enqueueer, containing its SPSC's minimum timestamp.

We apply some domain knowledge of Slotqueue algorithm to the definitions introduced in Section A.1.5.

Definition A.4.1.1 A **CAS-sequence** on a slot s of an enqueuee that affects s is the sequence of instructions from Line 12 to Line 17 of its `refreshEnqueue` (Procedure 22).

Definition A.4.1.2 A **slot-modification instruction** on a slot s of an enqueuee that affects s is Line 17 of `refreshEnqueue` (Procedure 22).

Definition A.4.1.3 A **CAS-sequence** on a slot s of a dequeuee that affects s is the sequence of instructions from Line 47 to Line 51 of its `refreshDequeue` (Procedure 25).

Definition A.4.1.4 A **slot-modification instruction** on a slot s of a dequeuee that affects s is Line 51 of `refreshDequeue` (Procedure 25).

Definition A.4.1.5 A **CAS-sequence** of a dequeuee/enqueuee is said to **observe a slot value of s_0** if it loads s_0 at Line 12 of `refreshEnqueue` or Line 47 of `refreshDequeue`.

The followings are some other definitions that will be used throughout our proof.

Definition A.4.1.6 For an enqueuee or dequeuee op , $rank(op)$ is the rank of the enqueueer whose local SPSC is affected by op .

Definition A.4.1.7 For an enqueueer whose rank is r , the value stored in its corresponding slot at time t is denoted as $slot(r, t)$.

Definition A.4.1.8 For an enqueueer with rank r , the minimum timestamp among the elements between First and Last in its local SPSC at time t is denoted as $min-spvc-ts(r, t)$.

Definition A.4.1.9 For an enqueue, **slot-refresh phase** refer to its execution of Line 3 - Line 4 of Procedure 21.

Definition A.4.1.10 For a dequeue, **slot-refresh phase** refer to its execution of Line 26 - Line 27 of Procedure 23.

Definition A.4.1.11 For a dequeue, **slot-scan phase** refer to its execution of Line 29 - Line 43 of Procedure 24.

A.4.2 Correctness

A.4.2.1 ABA problem

A.4.2.2 Memory reclamation

Notice that Slotqueue pushes the memory reclamation problem to the underlying SPSC. If the underlying SPSC is memory-safe, Slotqueue is also memory-safe.

A.4.2.3 Linearizability

Procedure 26: bool enqueue(data_t v)

```
1 timestamp = faa(Counter, 1)
2 if (!spvc_enqueue(&Spvc, (v, timestamp))) return false
3 if (!refreshEnqueue(timestamp))
4   | refreshEnqueue(timestamp)
5 return true
```

Procedure 27: bool refreshEnqueue(timestamp_t ts)

```
6 front = (data_t {}, timestamp_t {})  
7 success = spsc_readFront(Spsc, &front)  
8 new_timestamp = success ? front.timestamp : MAX_TIMESTAMP  
9 if (new_timestamp != ts)  
10 | return true  
11 old_timestamp = timestamp_t {}  
12 read(Slots + Self_rank, &old_timestamp)  
13 success = spsc_readFront(Spsc, &front)  
14 new_timestamp = success ? front.timestamp : MAX_TIMESTAMP  
15 if (new_timestamp != ts)  
16 | return true  
    return cas(Slots + Self_rank,  
17    old_timestamp,  
    new_timestamp)
```

Procedure 28: bool dequeue(data_t* output)

```
19 rank = readMinimumRank()  
20 if (rank == DUMMY_RANK)  
21 | return false  
22 output_with_timestamp = (data_t {}, timestamp_t {})  
23 if (!spsc_dequeue(Spsc, &output_with_timestamp))  
24 | return false  
25 *output = output_with_timestamp.data  
26 if (!refreshDequeue(rank))  
27 | refreshDequeue(rank)  
28 return true
```

Procedure 29: uint64_t readMinimumRank()

```
29 buffered_slots = timestamp_t[Process_count] {}
30 for index in 0..Process_count
31   | read(Slots + index, &buffered_slots[index])
32 if every entry in buffered_slots is MAX_TIMESTAMP
33   | return DUMMY_RANK
34 let rank be the index of the first slot that contains the minimum timestamp
   among buffered_slots
35 for index in 0..rank
36   | read(Slots + index, &buffered_slots[index])
37 min_timestamp = MAX_TIMESTAMP
38 for index in 0..rank
39   | timestamp = buffered_slots[index]
40   | if (min_timestamp < timestamp)
41     | min_rank = index
42     | min_timestamp = timestamp
43 return min_rank
```

Procedure 30: refreshDequeue(rank: int) **returns** bool

```
46 old_timestamp = timestamp_t {}
47 read(Slots + rank, &old_timestamp)
48 front = (data_t {}, timestamp_t {})
49 success = spsc_readFront(&Spscs[rank], &front)
50 new_timestamp = success ? front.timestamp : MAX_TIMESTAMP
   return cas(Slots + rank,
51   old_timestamp,
   new_timestamp)
```

A.4.3 Progress guarantee

Notice that every loop in Slotqueue is bounded, and no method have to wait for another. Therefore, Slotqueue is wait-free.

A.4.4 Theoretical performance

A summary of the theoretical performance of Slotqueue is provided in Table 7, which is already shown in Table 4. By R , we mean remote operations and by L we mean local operations.

Operations	Time-complexity
enqueue	$4R + 3L$
dequeue	$3R + 2nL$

Table 7: Theoretical performance summary of Slotqueue. R means remote operations and L means local operations.

For enqueue, we consider Procedure 21. Line 1 causes 1 truly remote operation, as the distributed counter is hosted on the dequeuer. Line 2, as discussed in the theoretical performance of SPSC, causes $R + L$ operations. In the worst case, two refreshEnqueue calls are executed. We then consider each refreshEnqueue call. Line 7 causes $R + L$ operations. Most of the time, Line 12 - Line 17 are not executed. Therefore, the two refreshEnqueue calls cause at most $2R$ operations. So in total, $4R + 3L$ operations are required.

For dequeue, we consider Procedure 23. Line 19 causes most of the remote operations: The double scan of the Slots array causes about $2nL$ operations. We consider the truly remote operations. Line 23 causes $R + L$ operations. The double retry on Line 26 - Line 27 each causes L operation (Line 47) and R operation. So in total, $3R + 2nL$ operations are required.

References

- [1] J. Schuchart, A. Bouteiller, and G. Bosilca, "Using MPI-3 RMA for Active Messages," 2019. doi: [10.1109/ExaMPI49596.2019.00011](https://doi.org/10.1109/ExaMPI49596.2019.00011).
- [2] Message Passing Interface Forum, *MPI: A Message-Passing Interface Standard Version 3.1*. 2015. [Online]. Available: <https://www.mpi-forum.org/docs/mpi-4.1/mpi41-report.pdf>
- [3] M. Baker and R. Buyya, "Cluster computing: the commodity supercomputer," 1999, *Wiley Online Library*.
- [4] S. Mullender, *Distributed systems*. Association for Computing Machinery, 1990.
- [5] M. Herlihy and N. Shavit, *The Art of Multiprocessor Programming, Revised Reprint*. Morgan Kaufmann, 2012.
- [6] M. Herlihy and N. Shavit, "On the Nature of Progress," 2011, *Springer Berlin Heidelberg*.
- [7] F. Ellen and T. Brown, "Concurrent Data Structures," 2016, *Association for Computing Machinery*. doi: [10.1145/2933057.2933123](https://doi.org/10.1145/2933057.2933123).
- [8] C. Hewitt, P. Bishop, and R. Steiger, "A universal modular ACTOR formalism for artificial intelligence," 1973, *Morgan Kaufmann Publishers Inc*.
- [9] A. Malhotra, "Concurrency Patterns in Golang: Real-World Use Cases and Performance Analysis," 2025.
- [10] Q. Yang, L. Tang, Y. Guo, N. Kuang, S. Zhong, and H. Luo, "WRLqueue: A Lock-Free Queue For Embedded Real-Time System," 2022. doi: [10.1109/HPCC-DSS-SmartCity-DependSys57074.2022.00197](https://doi.org/10.1109/HPCC-DSS-SmartCity-DependSys57074.2022.00197).
- [11] J. Wang, Q. Jin, X. Fu, Y. Li, and P. Shi, "Accelerating Wait-Free Algorithms: Pragmatic Solutions on Cache-Coherent Multicore Architectures," 2019. doi: [10.1109/ACCESS.2019.2920781](https://doi.org/10.1109/ACCESS.2019.2920781).
- [12] P. Jayanti and S. Petrovic, "Logarithmic-time single deleter, multiple inserter wait-free queues and stacks," 2005, *Springer-Verlag*. doi: [10.1007/11590156_33](https://doi.org/10.1007/11590156_33).
- [13] D. Adas and R. Friedman, "A Fast Wait-Free Multi-Producers Single-Consumer Queue," 2022, *Association for Computing Machinery*. doi: [10.1145/3491003.3491004](https://doi.org/10.1145/3491003.3491004).
- [14] B. Brock, A. Buluç, and K. Yelick, "BCL: A Cross-Platform Distributed Data Structures Library," 2019, *Association for Computing Machinery*. doi: [10.1145/3337821.3337912](https://doi.org/10.1145/3337821.3337912).
- [15] Thanh-Dang Diep, Phuong Hoai Ha, and Karl Furlinger, "A general approach for supporting nonblocking data structures on distributed-memory systems," 2023. doi: <https://doi.org/10.1016/j.jpdc.2022.11.006>.

- [16] H. Devarajan, A. Kougkas, K. Bateman, and X.-H. Sun, "HCL: Distributing Parallel Data Structures in Extreme Scales," 2020. doi: [10.1109/CLUSTER49012.2020.00035](https://doi.org/10.1109/CLUSTER49012.2020.00035).
- [17] G. Dewan and L. Jenkins, "Paving the way for Distributed Non-Blocking Algorithms and Data Structures in the Partitioned Global Address Space model," 2020. doi: [10.1109/ipdpsw50202.2020.00111](https://doi.org/10.1109/ipdpsw50202.2020.00111).
- [18] J. Feo, O. Villa, A. Tumeo, and S. Secchi, "Irregular applications: architectures & algorithms." Association for Computing Machinery, pp. 1–2, 2011. doi: [10.1145/2089142.2089144](https://doi.org/10.1145/2089142.2089144).
- [19] W. Gropp, R. Lusk, R. Ross, and R. Thakur, "Advanced MPI: I/O and one-sided communication." p. 202, 2006.
- [20] J. Dinan, P. Balaji, D. Buntinas, D. Goodell, W. Gropp, and R. Thakur, "An implementation and evaluation of the MPI 3.0 one-sided communication interface," 2016, *John Wiley and Sons Ltd*. doi: [10.1002/cpe.3758](https://doi.org/10.1002/cpe.3758).
- [21] M. P. Herlihy and J. M. Wing, "Linearizability: a correctness condition for concurrent objects," 1990, *Association for Computing Machinery*. doi: [10.1145/78969.78972](https://doi.org/10.1145/78969.78972).
- [22] M. Herlihy, "Wait-free synchronization," 1991, *Association for Computing Machinery*. doi: [10.1145/114005.102808](https://doi.org/10.1145/114005.102808).
- [23] M. M. Michael and M. L. Scott, "Simple, fast, and practical non-blocking and blocking concurrent queue algorithms." Association for Computing Machinery, pp. 267–275, 1996.
- [24] M. M. Michael, "Safe memory reclamation for dynamic lock-free objects using atomic reads and writes." Association for Computing Machinery, pp. 21–30, 2002.
- [25] J. V. Guttag, J. J. Horning, and J. M. Wing, "The Larch Family of Specification Languages," *IEEE Software*, vol. 2, no. 5, pp. 24–36, Sep. 1985.
- [26] M. P. Herlihy and J. M. Wing, "Axioms for concurrent objects." Association for Computing Machinery, pp. 13–26, 1987. [Online]. Available: <https://doi.org/10.1145/41625.41627>
- [27] S. Owicki and D. Gries, "An axiomatic proof technique for parallel programs I," *Acta Informatica*, vol. 6, pp. 319–340, 1976.