

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 DISTRIBUTED BARRIER
ALGORITHMS USING THE HYBRID PROGRAMMING
MODEL COMBINING MPI-3 AND C++11**

Major: Computer Science

THESIS COMMITTEE: 12
MEMBER SECRETARY: NGUYỄN QUANG HÙNG
SUPERVISORS: THOẠI NAM
DIỆP THANH ĐĂNG

—o0o—

STUDENT: PHẠM VÕ QUANG MINH -
2111762

HCMC, 01/2025

Disclaimers

I hereby declare that this specialized project is the result of my own research and experiments. It has not been copied from any other sources. All content presented and implemented within this document reflects my own hard work, dedication, and honesty, conducted under the 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.

All data, references, and sources have been legally cited and are explicitly mentioned in the footnotes and references section.

I accept full responsibility for the accuracy of the claims and content in this specialized project and am willing to face any consequences or penalties should any violations or misconduct be identified.

Acknowledgements

First, I want to express my gratitude to my school and teachers for providing me with a strong foundation in the field of computer science and software engineering. Their guidance has helped me refine my skills and develop my mindset as an engineer and scientist. These experiences have not only supported me in this specialized project but will also guide me throughout my lifelong journey.

I am especially thankful to my mother and brother for their unwavering support during both good times and bad. They have always been there for me, offering emotional and financial support and being the pillars I can rely on when times are tough.

I would also like to extend my deepest gratitude to Mr. Diệp Thanh Đăng for inspiring me and guiding me toward this research topic. His mentorship helped me discover my true passion and led me to choose this meaningful topic for my thesis.

Contents

Chapter I: Introduction	8
1.1 Motivation	9
1.2 Objectives	9
1.3 Scope	9
1.4 Project report structure	10
Chapter II: Background	12
2.1 Message Passing Interface (MPI)	13
2.1.1 Introduction	13
2.1.2 One-Sided Communications	13
2.1.3 Memory Model	14
2.2 C++11 Multithreading	15
2.2.1 Core Components	15
2.2.2 Thread Management	15
2.2.2.1 Creating Threads	16
2.2.2.2 Safe Thread Joining	16
2.2.3 Synchronization	17
2.2.3.1 Basic Mutex Operations	18
2.2.3.2 RAII Lock Guards	18
2.2.3.3 Deadlock Prevention	18
2.2.3.4 Condition Variables	19
2.2.4 Atomic Operations	19
2.2.5 Future and Promise	20
2.2.6 Key Advantages	21
2.3 HPC Context	21
Chapter III: Related Works	23
3.1 Barrier Synchronization Algorithm Selection	24
3.2 Hybrid Parallelization Approaches	24
Chapter IV: Adaptation from shared memory to distributed memory	25
4.1 Brook Algorithm	26
4.2 My proposed implementation of Brook's algorithm	26
4.2.1 Notes on the adapted algorithm	27
Chapter V: Preliminary Results	28
5.1 Test Environment	29
5.2 Test Implementation	29
5.3 Compilation and Execution	30
5.4 Results	31
Chapter VI: Conclusions and Future Works	32
6.1 Accomplishments This Semester	33
6.2 Challenges and Learnings	33



6.3 Planned Research Trajectory	33
References	35

List of Listings

Listing 1: Basic thread creation in C++11	16
Listing 2: RAII-based safe thread joining Section 2.2	17
Listing 3: Basic Mutex in C++11	18
Listing 4: RAII-based mutex locking	18
Listing 5: Deadlock prevention using <code>scoped_lock</code>	19
Listing 6: Producer-consumer pattern using condition variables	19
Listing 7: Atomic operations with memory ordering	20
Listing 8: Asynchronous computation with future and promise	21
Listing 9: Barrier signature	27
Listing 10: My adaptation of Brook's two-process barrier using MPI primitives	27
Listing 11: Ran the barrier synchronization four times to prove the algorithm is reusable.	
30	



List of Images

Figure 1: Test results on Apple M1	31
Figure 2: Gantt chart timeline.	34

Chapter I: Introduction

Chapter 1 introduces the project's research topic, focusing on distributed barrier synchronization algorithms in high-performance computing.

Section 1.1 explores the motivation behind the research, highlighting the growing importance of high-performance computing across scientific domains, including large language model training and complex simulations. Particularly, the section emphasizes the critical role of barrier algorithms in synchronizing computational processes across different nodes.

Section 1.2 details the project objectives, which encompass a comprehensive exploration of high-performance computing concepts, including MPI-3 and C++11 threading technologies. The research aims to investigate a hybrid programming model that combines the communication strengths of MPI with the synchronization capabilities of C++11. The objectives include surveying existing barrier synchronization algorithms and proposing methods for their deployment on current high-performance computing systems.

Section 1.3 narrows down the project's focus, the scope includes exploring key HPC concepts, studying MPI-3 and C++11 threading features, and implementing a simple barrier synchronization algorithm using the hybrid programming model. This aims to equip the author with new concepts before tackling the research objectives.

The project report structure, outlined in Section 1.4, provides a clear roadmap for the research. The document is organized into six chapters, progressing from foundational concepts to detailed algorithm proposals, experimental results, and future research directions. Each chapter builds upon the previous one, creating a comprehensive narrative of the research journey from introduction to conclusion.

1.1 Motivation

Applications of high-performance computing have gained immense popularity due to its applications in various scientific domains, from simulations of particle movement in physics [1], weather prediction [2], and it has also gained tremendous popularity with the booming of Large Language Models (LLMs), with the rising demands of distributed machine learning where the models are trained on enormous datasets.

These applications require tremendous effort in communication and synchronization. One of the basic building blocks that is widely used by these applications is the barrier algorithm, which aims to synchronize computation to a point to exchange data between compute nodes. Therefore, improving the performance of barrier algorithms can lead to significant performance improvements in these applications.

Barrier algorithms are built upon many programming models. Recently, a new model has gained a lot of attention: a hybrid programming model between MPI and C++11 [3]. Specifically, it is a crossover between MPI's power of communicating and synchronizing across different compute nodes, and C++'s power of communicating and synchronizing within one compute node. Quaranta and Maddegedara [3] have shown the potential of this programming model in terms of performance through a newly proposed barrier algorithm. However, they have ignored a large number of existing barrier algorithms that have been designed for similar programming models. [4]

Barrier algorithms can be divided into two categories: shared-memory and distributed-memory. Traditionally, shared-memory barrier algorithms could not be run on distributed-memory systems. However, with the introduction of MPI-2's One-Sided Communication [5], it is now possible to run shared-memory barrier algorithms on distributed-memory systems.

1.2 Objectives

The main goal of this thesis is to research, compare, and potentially, invent distributed barrier synchronization algorithms using the hybrid programming model. Specifically:

- Survey existing barrier synchronization algorithms in related programming models.
- Implements variants of barrier synchronization algorithms in the hybrid programming model.
- Compare the performance of the implemented barrier algorithms with existing algorithms.
- Propose methods to deploy selected barrier algorithms on the current HPC system.

1.3 Scope

The scope of this Specialization Project, the project aims to provide and familiarize the author with new concepts and programming paradigms. The details are as follows:

- Explore key concepts in high-performance computing (HPC).
- Study MPI-3 and its programming model.
- Investigate C++11 threading and concurrency features.
- Research the hybrid programming model combining MPI-3 and C++11.
- Survey existing barrier synchronization algorithms.
- Implement a simple barrier synchronization algorithm with the hybrid model.

1.4 Project report structure

The rest of this report is organized as follows:

Chapter II: Background

This chapter establishes the theoretical and technical foundations of the research through three main sections:

Section 2.1 Message Passing Interface (MPI-3)

- Evolution and significance in distributed computing
- Advanced communication mechanisms, focusing on One-Sided Communications
- Remote Memory Access (RMA) operations and memory models
- Collective operations and parallel I/O capabilities

Section 2.2 Modern C++11 Multithreading

- Native threading support and platform independence
- Thread management and synchronization primitives
- Memory model and atomic operations
- Asynchronous programming features and their implications

Section 2.3 High-Performance Computing Context

- Integration of shared and distributed memory paradigms
- Hybrid programming models in modern HPC

Chapter III: Related Works

Surveys existing research and approaches to barrier synchronization algorithms.

Chapter IV: Adaptation from shared memory to distributed memory

This chapter introduces an adaptation of Brook's barrier algorithm, originally designed for shared memory models, to distributed memory systems using MPI's Remote Memory Access (RMA) operations.

This implementation leverages MPI RMA primitives for efficient communication in distributed memory systems, preserving the core synchronization logic of Brook's original algorithm.

Chapter V: Preliminary Results

Presents the implementation outcomes, experimental results, and performance analysis.

Chapter VI: Conclusions and Future Works

Summarizes accomplishments, outlines future research directions, and provides a timeline for planned activities.

Chapter II: Background

Chapter 2 provides a comprehensive background on the fundamental technologies underlying the research: Message Passing Interface (MPI) and C++11 Multithreading. The chapter explores these technologies within the context of high-performance computing, laying the groundwork for understanding the hybrid programming model proposed in the thesis.

Section 2.1 delves into MPI, presenting it as a critical message-passing library interface specification. The section goes beyond traditional communication models by highlighting MPI's advanced capabilities, including collective operations, remote-memory access, and parallel I/O. A focus is placed on One-Sided Communications, a communication mechanism that allows a single MPI process to independently manage communication tasks. The section elaborates on the various Remote Memory Access (RMA) operations and explores two distinct memory models: Separate and Unified Memory Models.

Section 2.2 explores C++11 Multithreading, marking a pivotal advancement in parallel computing. The section demonstrates how C++11 introduced native, platform-independent threading support through its standard library. It comprehensively covers key multithreading features, including thread management, synchronization primitives, atomic operations, and asynchronous programming constructs. The discussion emphasizes the advantages of C++11's threading model, such as type-safe synchronization and reduced dependency on platform-specific libraries.

The concluding Section 2.3 bridges these technologies within the high-performance computing context. It illustrates how C++11 multithreading and MPI can complement each other: C++11 provides efficient intra-node communication within a single node, while MPI handles inter-node communication between different computers in a cluster. This perspective sets the stage for the thesis's exploration of a hybrid programming model that leverages the strengths of both technologies.

By providing this detailed technological background, Chapter 2 establishes the theoretical and practical foundations necessary for understanding the subsequent research on barrier synchronization algorithms in high-performance computing.

2.1 Message Passing Interface (MPI)

This section explores Message Passing Interface, its various paradigms for message passing (Section 2.1.1). The section will then shift its focus into MPI's One-Sided Communication mechanics.

2.1.1 Introduction

MPI [6] is a message-passing library interface specification - essentially a set of standard guidelines for both implementors and users of the parallel programming model. In this programming paradigm, data traverses between the address spaces of different processes through cooperative operations - what we might call as “classical message-passing techniques”.

Beyond the classical model, MPI extends its capabilities by offering:

- Collective operations
- Remote-memory access operations
- Dynamic process creation
- Parallel I/O

Within this thesis, I will specifically dive deep into remote-memory access operations in MPI, which form the critical building block for MPI's One-Sided Communication, a concept first introduced in MPI-2 [5], and later refined in MPI-3 [7].

2.1.2 One-Sided Communications

Unlike “classical” Point-to-Point communication - where communication is a two-way collaborative effort between processes, with one sending and another receiving - Remote Memory Access introduces a more flexible communication mechanism called One-Sided communication.

In this approach, a single MPI process can now independently orchestrate the entire communication task, by specifying communication parameters for both the sending and receiving sides.

Regular send/receive communication requires matching operations by sender and receiver. To issue the matching operations, the application needs to distribute the transfer parameters. This may require all MPI processes to participate in a global computation, or to explicitly poll for potential communication requests to respond periodically. This could potentially introduce plenty overhead.

The use of RMA communication operations avoids the need for global computation or explicit polling [6].

Message-passing in general achieves two effects: communication of data from sender to receiver and synchronization of sender with receiver. RMA design separates these two functions.

Within MPI standards, the primitive communication operations are as follows:

- Remote write operations: MPI_PUT, MPI_RPUT
- Remote read operations: MPI_GET, MPI_RGET
- Remote update operations: MPI_ACCUMULATE, MPI_RACCUMULATE
- Combined read and update operations: MPI_GET_ACCUMULATE, MPI_RGET_ACCUMULATE, and MPI_FETCH_AND_OP
- Remote atomic swap: MPI_COMPARE_AND_SWAP

The primitive synchronization operations are:

- Active Target Communication: MPI_WIN_FENCE, MPI_WIN_POST, MPI_WIN_START, MPI_WIN_COMPLETE, MPI_WIN_WAIT
- Passive Target Communication: MPI_WIN_LOCK, MPI_WIN_UNLOCK

Example [6]: Implementing a critical region between multiple MPI processes with compare and swap. The call to MPI_WIN_SYNC is necessary on Process A after local initialization of A to guarantee the public copy has been updated with the initialization value found in the private copy. It would also be valid to call MPI_ACCUMULATE with MPI_REPLACE to directly initialize the public copy. A call to MPI_WIN_FLUSH would be necessary to assure A in the public copy of Process A had been updated before the barrier.

Process A

```
MPI_Win_lock_all
atomic location A
A=0
MPI_Win_sync
MPI_Barrier
stack variable r = 1
while(r != 0) do
r = MPI_Compare_and_swap(A, 0, 1)
MPI_Win_flush(A)
done
// critical region
r = MPI_Compare_and_swap(A, 1, 0)
MPI_Win_unlock_all
```

Process B

```
MPI_Win_lock_all

MPI_Barrier
stack variable r = 1
while(r != 0) do
r = MPI_Compare_and_swap(A, 0, 1)
MPI_Win_flush(A)
done
// critical region
r = MPI_Compare_and_swap(A, 1, 0)
MPI_Win_unlock_all
```

2.1.3 Memory Model

MPI supports two distinct memory models:

1. Separate Memory Model:

- No inherent memory consistency assumptions
- Similar to weakly coherent memory systems

- Requires explicit synchronization for correct memory access ordering

2. Unified Memory Model:

- Exploits cache-coherent hardware
- Supports hardware-accelerated one-sided operations
- Typically found in high-performance computing environments

The MPI design's flexibility allows implementors to leverage platform-specific communication mechanisms, including:

- Coherent and non-coherent shared memory
- Direct Memory Access (DMA) engines
- Hardware-supported put/get operations
- Communication coprocessors

2.2 C++11 Multithreading

C++11 (ISO/IEC 14882:2011) introduced native support for multithreading directly in the Standard Template Library (STL) [8]. This eliminated the previous reliance on platform-specific threading libraries like POSIX thread (pthread) and the Microsoft Windows API, offering a more portable and standardized approach to parallel computing. The introduction of these features marked a significant advancement in C++ development, providing developers with powerful, standardized tools for concurrent and parallel programming.

2.2.1 Core Components

The STL threading components include:

- Thread management (`<thread>`)
- Mutual exclusion (`<mutex>`)
- Condition variables (`<condition_variable>`)
- Atomic operations (`<atomic>`)
- Futures and promises (`<future>`)

2.2.2 Thread Management

The `<thread>` header has now become the pivot of C++11's multithreading support. It provides a lightweight, platform-independent mechanism for creating and managing threads

2.2.2.1 Creating Threads

```
#include <thread>
#include <iostream>

void worker_function() {
    std::cout << "Thread is running" << std::endl;
}

int main() {
    std::thread t(worker_function); // Create a thread
    t.join(); // Wait for the thread to complete
    return 0;
}
```

Listing 1: Basic thread creation in C++11

Important characteristics of `std::thread`:

- Non-copyable but movable
- Must be either joined or detached before destruction, or else `std::terminate()` will be called, resulting in immediate termination of the entire application. This is a safety feature to prevent accidental thread abandonment.
- Supports various callable objects (functions, lambdas, function objects)

2.2.2.2 Safe Thread Joining

When exceptions occur in the parent thread, proper thread management becomes crucial. To prevent unexpected program termination, we can use RAII patterns or C++20's `std::jthread`.


```
#include <thread>
#include <stdexcept>

class thread_guard {
    std::thread& t;
public:
    explicit thread_guard(std::thread& t_) : t(t_) {}
    ~thread_guard() {
        if(t.joinable()) {
            t.join();
        }
    }
    thread_guard(const thread_guard&) = delete;
    thread_guard& operator=(const thread_guard&) = delete;
};

void may_throw() {
    throw std::runtime_error("Exception occurred");
}

void worker_function() {
    // Thread work here
}

int main() {
    std::thread t(worker_function);
    thread_guard g(t); // RAII wrapper ensures join
    try {
        may_throw();
    }
    catch (...) {
        // Thread will be automatically joined
        throw;
    }
    return 0;
}
```

Listing 2: RAII-based safe thread joining Section 2.2

2.2.3 Synchronization

C++11 provides several RAII-compliant mechanisms to handle mutual exclusion and prevent race conditions.

2.2.3.1 Basic Mutex Operations

Mutexes (mutual exclusion objects) provide the most fundamental way to protect shared data. While direct lock/unlock operations are available, they should be used with caution as they don't provide automatic cleanup in case of exceptions.

```
std::mutex mtx;  
mtx.lock();    // Acquire lock  
// Critical section  
mtx.unlock();  // Release lock
```

Listing 3: Basic Mutex in C++11

2.2.3.2 RAII Lock Guards

RAII (Resource Acquisition Is Initialization) lock guards provide a safer alternative to manual mutex locking and unlocking. They automatically release the mutex when going out of scope, preventing resource leaks and deadlocks.

```
#include <mutex>  
  
class shared_resource {  
    std::mutex mtx;  
    int data;  
public:  
    void safe_operation() {  
        std::lock_guard<std::mutex> lock(mtx);  
        // Protected operations here  
    } // Automatically unlocked  
};
```

Listing 4: RAII-based mutex locking

2.2.3.3 Deadlock Prevention

Deadlocks occur when multiple threads are waiting for each other to release resources. C++11 provides `std::lock` to acquire multiple mutexes simultaneously using a deadlock avoidance algorithm. Later, C++17 introduced `std::scoped_lock` which combines this functionality with RAII principles, making it safer and more convenient to acquire multiple mutexes simultaneously while preventing deadlocks.

```
#include <mutex>

class complex_resource {
    std::mutex mtx1, mtx2;
public:
    void safe_operation() {
        std::scoped_lock lock(mtx1, mtx2); // Deadlock-free
        // Protected operations here
    }
};
```

Listing 5: Deadlock prevention using `scoped_lock`

2.2.3.4 Condition Variables

Condition variables enable threads to coordinate based on actual values or conditions rather than simply taking turns.

```
#include <condition_variable>
#include <mutex>
#include <queue>

std::queue<int> data_queue;
std::mutex mtx;
std::condition_variable cv;

void producer() {
    std::lock_guard<std::mutex> lock(mtx);
    data_queue.push(42);
    cv.notify_one();
}

void consumer() {
    std::unique_lock<std::mutex> lock(mtx);
    cv.wait(lock, []{ return !data_queue.empty(); });
    int data = data_queue.front();
    data_queue.pop();
}
```

Listing 6: Producer-consumer pattern using condition variables

2.2.4 Atomic Operations

Atomic operations provide lock-free programming capabilities, ensuring that certain operations are performed as a single, uninterruptible unit.

```
#include <atomic>

struct atomic_counter {
    std::atomic<int> value{0};

    void increment() {
        value.fetch_add(1, std::memory_order_relaxed);
    }

    bool compare_exchange(int expected, int desired) {
        return value.compare_exchange_strong(expected, desired,
                                              std::memory_order_acq_rel);
    }

    int load() const {
        return value.load(std::memory_order_acquire);
    }
};
```

Listing 7: Atomic operations with memory ordering

2.2.5 Future and Promise

C++11 provides powerful mechanisms for asynchronous programming:

```
#include <future>
#include <stdexcept>

std::promise<int> compute_promise;

void compute() {
    try {
        int result = /* complex computation */42;
        compute_promise.set_value(result);
    } catch (...) {
        compute_promise.set_exception(std::current_exception());
    }
}

int main() {
    std::future<int> result = compute_promise.get_future();
    std::thread t(compute);

    try {
        int value = result.get(); // Will throw if computation failed
        std::cout << "Result: " << value << std::endl;
    } catch (const std::exception& e) {
        std::cerr << "Computation failed: " << e.what() << std::endl;
    }

    t.join();
    return 0;
}
```

Listing 8: Asynchronous computation with future and promise

2.2.6 Key Advantages

- Platform-independent threading
- Standard library support
- Type-safe synchronization
- Reduced dependency on platform-specific libraries
- Simplified concurrent programming model

The introduction of these features provided developers with powerful, standardized API for concurrent and parallel programming without resorting to platform-specific libraries like pthread.

2.3 HPC Context

In the landscape of high-performance computing (HPC), this hybrid approach to use MPI for inter-node communication and multithreading for intra-node communication is not

new. However, third-party libraries like OpenMP and pthread have been said to introduced some hidden bugs and other issues regarding performance [3]. The introduction of C++11, which provides a standardized approach to shared-memory parallel processing, threading support has made it easier to write portable, reliable and efficient code.

Chapter III: Related Works

This section examines two key areas relevant to my implementation: barrier synchronization algorithms and hybrid parallelization approaches.

Section 3.1 review MellorCrummey's foundational work on barrier algorithms, which provides crucial insights into synchronization strategies for different multiprocessor architectures.

Section 3.2 explore recent advances in hybrid parallelization, particularly Quaranta and Maddegedara's novel approach combining MPI-3 shared memory windows with the C++11 memory model.

3.1 Barrier Synchronization Algorithm Selection

MellorCrummey et al [9] proposed nuanced recommendations for barrier synchronization:

1. Broadcast-Based Cache-Coherent Multiprocessors:
 - For modest processor counts: Utilize a centralized counter-based barrier
 - For larger scales: Implement a 4-ary arrival tree with a central sense-reversing wakeup flag
2. Multiprocessors Without Coherent Caches or With Directory-Based Coherency:
 - Dissemination Barrier [10]:
 - Distributed data structures respecting locality
 - Critical path approximately one-third shorter than tree-based barriers
 - Total interconnect traffic complexity of $O(P \log P)$
 - Tree-Based Barrier [9]:
 - Alternative approach with different performance characteristics
 - Total interconnect traffic complexity of $O(P)$

The dissemination barrier demonstrates superior performance on architectures like the Butterfly, which can execute parallel non-interfering network transactions across multiple processors.

3.2 Hybrid Parallelization Approaches

Quaranta and Maddegedara [3] proposed a novel hybrid approach combining MPI-3 shared memory windows with C11/C++11 memory model. Their work demonstrates several key advantages:

- Efficient intranode communication using MPI shared memory windows
- Fine-grained synchronization using C++11 atomic operations
- Reduced variance in execution times
- Enhanced synchronization between processes, especially in multi-node environments
- Significant performance improvements in ghost updates compared to flat MPI
- More efficient synchronization of shared data compared to RMA-based methods

Chapter IV: Adaptation from shared memory to distributed memory

This chapter presents my proposed adaptation of Brook's barrier algorithm within the context of distributed memory models.

Section 4.1 presents Brook's barrier algorithm within the context of the shared memory model.

Section 4.2 provides a straightforward implementation of Brook's two-process barrier algorithm in a distributed memory programming model using MPI's Remote Memory Access (RMA) operations.

4.1 Brook Algorithm

Brook [11] bases the n-process barrier on a two-process barrier using two shared variables. The algorithm is as follows:

Brook-Barrier-Algorithm

```
1 procedure Brook's-Barrier(SetByMyProcess, SetByTargetProcess)
2   while SetByMyProcess is true do wait
3   SetByMyProcess  $\leftarrow$  true
4   while SetByTargetProcess is false do wait
5   SetByTargetProcess  $\leftarrow$  false
6 end procedure
```

We can visualize the algorithm as follows:

Step	Process 1	Process 2
1	while SetByProcess1 do wait;	while SetByProcess2 do wait;
2	SetByProcess1 := true;	SetByProcess2 := true;
3	while not SetByProcess2 do wait;	while not SetByProcess1 do wait;
4	SetByProcess2 := false;	SetByProcess1 := false;

4.2 My proposed implementation of Brook's algorithm

We can extend the concept of Brooks' two-process barrier, where two processes share memory through shared variables, to the one-sided communication model.

In this model, one process can directly access the variables of another process. Instead of storing the shared variables in a common memory segment, each process maintains its own copy, allowing the other process to read and modify it using one-sided communication primitives.

This is my simple adaptation of Brook's two-process barrier algorithm to a distributed memory model using MPI RMA operations:

```
MPI_BROOK_BARRIER(exposed_flag, win, target_rank)
```

INPUTS:

exposed_flag: BOOL

win: MPI_WIN

target_rank: INT

OUTPUTS:

None

Listing 9: Barrier signature

Distributed-Mem-Brook-Barrier

```
1 procedure Brook's-Barrier(exposed_flag, win, target_rank)
2   while exposed_flag is true do wait
3   exposed_flag ← true
4   while target_flag is false do wait
5     MPI_WIN_LOCK(win)
6     MPI_GET_ACCUMULATE(&target_flag, 0, BOOL, &target_flag, 1, BOOL,
7       target_rank, 0, 1, BOOL, MPI_NO_OP, win);
8     MPI_WIN_FLUSH(win)
9     MPI_WIN_UNLOCK(win)
10  end while
11  false_value ← false
12  MPI_WIN_LOCK(win)
13  MPI_ACCUMULATE(&>false_value, 1, BOOL, target_rank, 0, 1, BOOL,
14    MPI_REPLACE, win);
15  MPI_WIN_FLUSH(win)
16  MPI_WIN_UNLOCK(win)
17 end procedure
```

Listing 10: My adaptation of Brook's two-process barrier using MPI primitives

4.2.1 Notes on the adapted algorithm

The implementation uses `MPI_GET_ACCUMULATE` and `MPI_ACCUMULATE` instead of `MPI_GET` and `MPI_PUT` for atomicity guarantees. The `MPI_WIN_LOCK` and `MPI_WIN_UNLOCK` functions are used to ensure atomicity of the operations.

Chapter V: Preliminary Results

This chapter showcases the preliminary test result of the Brook's two-process barrier algorithm using MPI's RMA operations.

5.1 Test Environment

- Hardware: Apple M1 chip
- Compiler: MPICH's mpic++ compiler
- MPI Implementation: MPICH
- Compilation flags: -std=c++11
- Operating System: macOS

5.2 Test Implementation

The test program was designed to verify the basic functionality of the barrier synchronization:

```
int main(int argc, char **argv) {
    // init the mpi world
    MPI_Init(&argc, &argv);
    MPI_Comm comm = MPI_COMM_WORLD;

    // get number of ranks
    int n_ranks;
    MPI_Comm_size(comm, &n_ranks);
    // get current rank
    int rank;
    MPI_Comm_rank(comm, &rank);

    if (n_ranks != 2) {
        if (rank == 0) {
            fprintf(stderr, "This program requires exactly 2 processes.\n");
        }
        MPI_Abort(comm, 1);
    }

    // set the target
    int target_rank = 1 - rank;

    // create a window
    bool exposed_bool{false};
    MPI_Win win_buffer_handler;
    MPI_Win_create(&exposed_bool, sizeof(bool), sizeof(bool),
        MPI_INFO_NULL, comm,
        &win_buffer_handler);

    barrier_brook(exposed_bool, win_buffer_handler, target_rank);
    barrier_brook(exposed_bool, win_buffer_handler, target_rank);
    barrier_brook(exposed_bool, win_buffer_handler, target_rank);
    barrier_brook(exposed_bool, win_buffer_handler, target_rank);

    printf("Process %d: reached destination\n", rank);

    return MPI_Finalize();
}
```

Listing 11: Ran the barrier synchronization four times to prove the algorithm is reusable.

5.3 Compilation and Execution

The program was compiled using the following command:

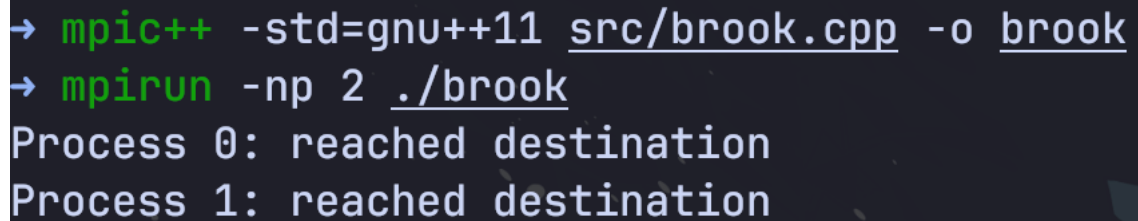
```
mpic++ -std=c++11 src/brook.cpp -o brook
```

Execution was performed using MPICH's mpirun with two processes:

```
mpirun -np 2 ./brook
```

5.4 Results

The test results demonstrated successful barrier synchronization between two processes:



```
→ mpic++ -std=gnu++11 src/brook.cpp -o brook
→ mpirun -np 2 ./brook
Process 0: reached destination
Process 1: reached destination
```

Figure 1: Test results on Apple M1

Chapter VI: Conclusions and Future Works

This chapter presents the accomplishments and future direction of the research.

Section 6.1 highlights the progress made this semester, including the implementation of Brook's two-process barrier algorithm using MPI's RMA operations.

Section 6.2 discusses the challenges faced and the learnings acquired during the semester.

Section 6.3 outlines the planned research trajectory, focusing on further algorithm development, benchmarking, and integration into larger systems. This chapter concludes with a Gantt chart timeline illustrating the key milestones and activities for the upcoming semester.

6.1 Accomplishments This Semester

Within the confines of this semester, I have:

- Studied and familiarized myself with MPI and its One-Sided Communication techniques
- Explored C++11 threading and concurrency features
- Successfully implemented the two-process Brooks barrier algorithm utilizing MPI's Remote Memory Access (RMA) Operations

6.2 Challenges and Learnings

My journey into parallel programming has been both challenging and enlightening. Initially unfamiliar with the concepts, I faced a steep learning curve, particularly when delving into the MPI standard documentation and its historical context. Understanding the intricacies of memory models and synchronization techniques required considerable effort, including extensive reading and iterative comprehension of the MPI standard.

Grasping the underlying memory model and the mechanics of one-sided communication was crucial. Although the adaptation itself was straightforward, the preparatory work - building a robust conceptual framework - posed a significant challenge.

Despite the challenges, this process has been highly rewarding. I have developed a solid foundation in distributed memory models and synchronization techniques. However, my work remains a stepping stone; while I successfully adapted Brook's barrier algorithm to distributed memory, I have yet to achieve a fully functional implementation for hybrid memory models. This leaves room for future exploration and growth in this domain.

6.3 Planned Research Trajectory

For the upcoming semester, my research will focus on:

- Understanding of C++ and MPI Memory model to implement C++ and MPI Hybrid model
- Adapting and Implementing barrier algorithms using MPI's One-Sided Communication
- Testing of proposed algorithms on existing clusters
- Performance benchmarking to identify and select optimal barrier synchronization strategies
- Integration and contribution to existing applications

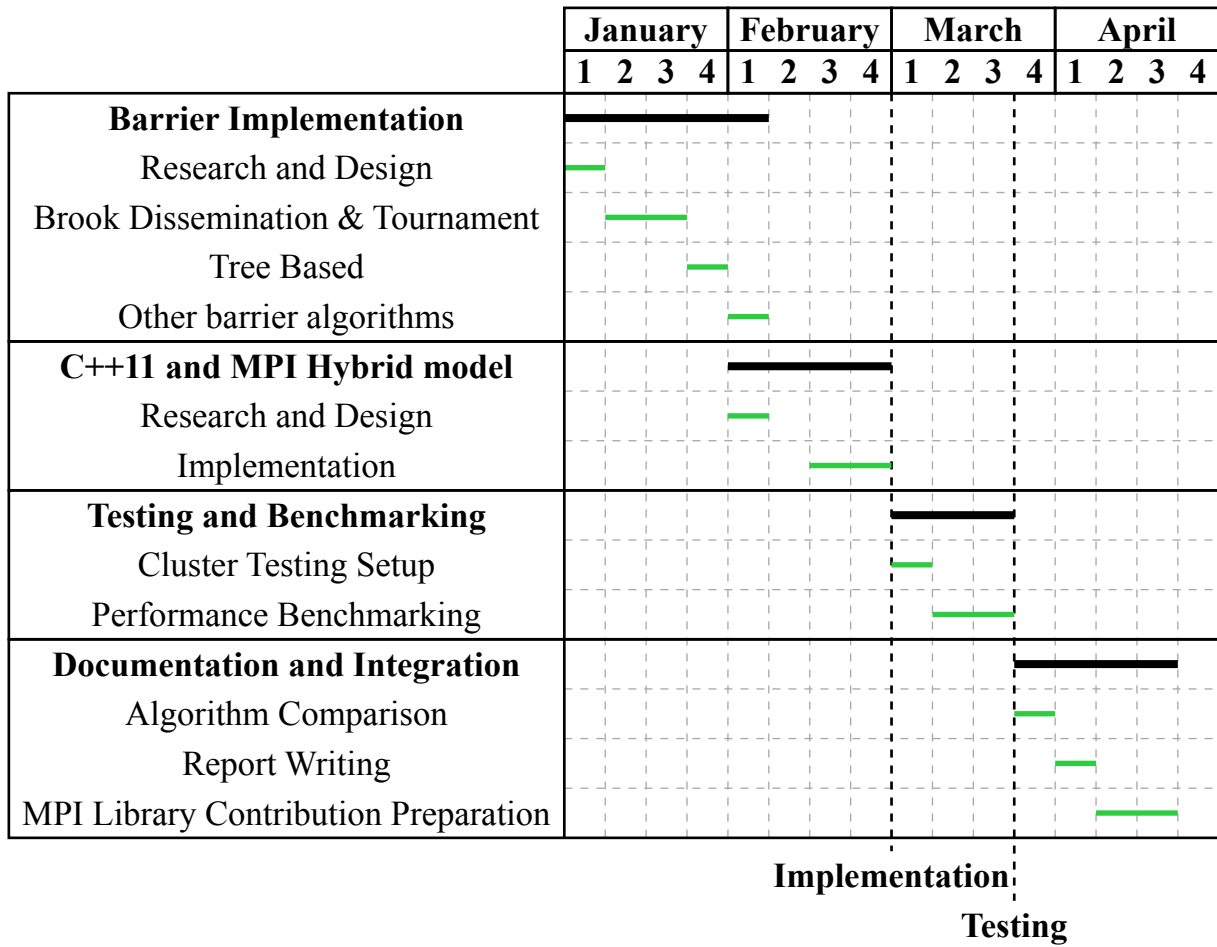


Figure 2: Gantt chart timeline.

References

- [1] “High Performance Computing (HPC) Applications and Examples.” Accessed: Nov. 26, 2024. [Online]. Available: <https://phoenixnap.com/kb/hpc-applications>
- [2] “NASA Global Weather Forecasting.” Accessed: Jan. 01, 2025. [Online]. Available: <https://www.nccs.nasa.gov/sci-tech/case-studies/nasa-global-weather-forecasting>
- [3] L. Quaranta and L. Maddegedara, “A Novel MPI+MPI Hybrid Approach Combining MPI-3 Shared Memory Windows and C11/C++11 Memory Model,” *Journal of Parallel and Distributed Computing*, vol. 157, pp. 125–144, Nov. 2021, doi: [10.1016/j.jpdc.2021.06.008](https://doi.org/10.1016/j.jpdc.2021.06.008).
- [4] M. Herlihy, N. Shavit, V. Luchangco, and M. Spear, *The Art of Multiprocessor Programming*, 2nd ed. Morgan Kaufman, 2020.
- [5] Message Passing Interface Forum, *MPI: A Message-Passing Interface Standard, Version 2.1*. High Performance Computing Center Stuttgart (HLRS), 2008.
- [6] Message Passing Interface Forum, *MPI: A Message-Passing Interface Standard Version 4.1*. Accessed: Nov. 2023. [Online]. Available: <https://www.mpi-forum.org/docs/mpi-4.1/mpi41-report.pdf>
- [7] Message Passing Interface Forum, *MPI: A Message-Passing Interface Standard Version 3.0*. High Performance Computing Center Stuttgart (HLRS), 2012.
- [8] A. Williams, *C++ Concurrency in Action*, 2nd ed. Manning, 2019, pp. 1–172.
- [9] “Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors,” *Association for Computing Machinery*. doi: [10.1145/103727.103729](https://doi.org/10.1145/103727.103729).
- [10] *An Optimal Scheme for Disseminating Information*. University of Kentucky, Department of Computer Science, 1988.
- [11] E. D. Brooks, “The butterfly barrier,” *International Journal of Parallel Programming*, vol. 15, no. 4, pp. 295–307, Aug. 1986, doi: [10.1007/BF01407877](https://doi.org/10.1007/BF01407877).