

Group3@IAU-ARTI503:

Parallel Text Processing and Word Frequency Counter

Abdullah Albattat¹, Hussain Alghubari¹, Muhannad Almahmoud¹,
Khalid Alghamdi¹, and Abdullah Aladwani¹

¹Imam Abdulrahman Bin Faisal University
College of Computer Science and Information Technology
Dammam, Saudi Arabia

ARTI503 - Parallel Computer Architecture
December 10, 2025

Contents

1	Introduction	3
1.1	Objective and Purpose	3
1.1.1	Problem Statement	3
1.1.2	Project Purpose	4
1.1.3	Why This Problem Matters for Parallel Computing	4
1.2	Code Selection and Justification	5
1.2.1	Why We Selected This Problem	5
1.2.2	Suitability for OpenMP Parallelization	7
1.2.3	Computational Complexity Analysis	7
1.2.4	Code Origin	8
2	Sequential Benchmarking	9
2.1	Benchmarking Methodology	9
2.1.1	Measurement Tools and Techniques	9
2.1.2	Test Environment Specifications	10
2.1.3	Dataset Specifications	10
2.2	Performance Results	11
2.2.1	Execution Time Analysis	11
2.2.2	Throughput Analysis	11
2.2.3	Performance Visualization	12
2.3	Bottleneck Identification	14
2.3.1	Profiling Analysis	14
2.3.2	Execution Time Breakdown	17
2.4	Scalability Analysis	18
2.4.1	Linear Scaling Verification	18
2.4.2	Performance Baseline Summary	18
3	Parallel Implementation	19
3.1	Parallelization Strategy and Methodology	19
3.1.1	OpenMP Framework Selection	19
3.1.2	Parallel Algorithm Design	19
3.1.3	Key Parallelization Techniques	21
3.2	Synchronization Methods Implemented	22
3.2.1	Method 1: Reduction (Default)	22
3.2.2	Method 2: Atomic Operations	22
3.2.3	Method 3: Critical Sections	23
3.3	Race Condition Analysis	23
3.3.1	Identified Race Conditions	23
3.3.2	Synchronization Method Comparison	24
4	Parallel Performance Results	25
4.1	Benchmarking Methodology	25
4.1.1	Experimental Setup	25
4.1.2	Performance Metrics	25
4.2	Speedup and Efficiency Analysis	26
4.2.1	Best Performance Results	26
4.2.2	Speedup vs Thread Count	26
4.2.3	Efficiency Analysis	27
4.3	Synchronization Method Comparison	28
4.4	Scalability Discussion	29
4.4.1	Strong Scaling Analysis	29
4.4.2	Bottleneck Analysis	30
5	Discussion and Limitations	31
5.1	Performance Analysis	31
5.1.1	Why Speedups Are Modest	31
5.1.2	Comparison with Expectations	31
5.2	Challenges and Debugging	32
5.2.1	Race Conditions Encountered	32
5.2.2	Debugging Techniques Used	32
5.3	Limitations of Current Approach	33
5.3.1	Technical Limitations	33
5.3.2	Algorithmic Limitations	33
5.4	Future Work and Improvements	34

5.4.1	Short-Term Improvements	34
5.4.2	Long-Term Research Directions	35
6	Template Usage and Documentation	36
6.1	Template Completion	36
6.1.1	Project Completion Status	36
6.2	Code Visualization and Repository	38
6.2.1	Project Structure	38
6.2.2	Code Access Methods	38
6.2.3	Code Functionality Explanation	39
6.3	Document Formatting and Compilation	40
6.3.1	LaTeX Formatting and Compilation	40
6.4	Project Implementation Summary	40
6.4.1	Parallel Implementation Approach	40
7	References	41
8	Appendix	42
8.1	Appendix A: Project Repository Structure	42
8.2	Appendix B: Build and Run Instructions	43
8.3	Appendix C: Performance Metrics Definitions	44
8.4	Appendix D: Synchronization Method Code Examples	45
8.5	Appendix E: Acknowledgments	46

1 Introduction

1.1 Objective and Purpose

1.1.1 Problem Statement

In the era of big data and digitalization, text processing has become a fundamental computational task across many domains like:

- **Search Engines:** Indexing billions of web pages
- **Social Media Analytics:** Processing real-time streams of user-generated content
- **Natural Language Processing:** Processing large corpora for machine learning models
- **Document Management Systems:** Managing and analyzing enterprise documents
- **Log Analysis:** Processing server and application logs for monitoring

Conventional sequential text-processing algorithms are severely limited for processing large-scale data. As file sizes increase from megabytes to gigabytes, single-threaded methods turn into major bottlenecks, leading to:

- Unacceptable processing latency
- Underutilization of modern multi-core processors
- Poor scalability for growing data volumes
- Inability to meet real-time processing requirements

1.1.2 Project Purpose

The primary purpose of this project is to **design, implement, and evaluate** a parallel word frequency counter that leverages multi-core processors to dramatically improve the performance of text processing.

Specific Objectives:

1. **Implement Sequential Baseline:** Create a sequential word counter to establish performance baselines
2. **Develop Parallel Solution:** Develop an OpenMP-based parallel solution through shared-memory programming
3. **Performance Analysis:** Measure and compare execution time, speedup, and efficiency across different configurations
4. **Scalability Study:** Evaluate how performance scales with:
 - Dataset size (10 MB to 100+ MB)
 - Thread count (2, 4, 8, 16 threads)
5. **Optimization:** Identify and minimize synchronization overhead and load imbalance
6. **Real-World Application:** Demonstrate practical benefits of parallelization in text analytics

1.1.3 Why This Problem Matters for Parallel Computing

Word frequency counting is an ideal candidate for parallel computing education and research because:

- **Computational Intensity:** Processing large text files involves millions of operations
- **Data Independence:** Text chunks can be processed independently
- **Clear Metrics:** Speedup and efficiency are easily measurable
- **Real-World Relevance:** Directly applicable to industry problems
- **Educational Value:** Demonstrates key parallel computing concepts (synchronization, load balancing, Amdahl's Law)

1.2 Code Selection and Justification

1.2.1 Why We Selected This Problem

We selected the word frequency counter problem for parallelization because it exhibits **ideal characteristics** for demonstrating parallel computing principles:

1. Computational Intensity:

- Processing large text files (50-100 MB) involves **tens of millions of operations**
- Each word requires: reading, tokenization, normalization, and hash map insertion
- Sequential processing can take several seconds to minutes for large datasets

2. Existence of Parallelizable Loops:

The core algorithm contains **highly parallelizable loops**:

```
1 // Main processing loop - SEQUENTIAL BOTTLENECK
2 while (file >> word) {
3     // Step 1: Normalize word (lowercase, remove punctuation)
4     normalized_word = normalize(word); // CPU-intensive
5
6     // Step 2: Update frequency map
7     wordFreq[normalized_word]++;      // Hash map operation
8
9     totalWords++;
10 }
```

Listing 1: Sequential Word Counting Loop - Main Bottleneck

This loop is executed **millions of times** and dominates execution time, making it the prime target for parallelization.

3. Data Independence:

- Text can be divided into independent chunks
- Each chunk can be processed in parallel without dependencies
- Only the final merge step requires synchronization

4. Multiple Parallelization Opportunities:

Table 1: Parallelization Opportunities in Word Frequency Counter

Component	Sequential Operation	Parallel Strategy
File Reading	Read entire file sequentially	Divide into chunks, parallel read
Tokenization	Process words one-by-one	Each thread processes its chunk
Word Normalization	Sequential character processing	Independent per word
Frequency Counting	Single global hash map	Thread-local maps + merge
Result Sorting	Sequential sort	Parallel sort (if needed)

5. Presence of Different Data Types:

The code involves multiple data types suitable for parallel processing:

- `std::string`: Text data requiring character-level operations
- `std::unordered_map<string, unsigned long long>`: Hash table for frequency storage
- `unsigned long long`: Large integer counters
- `std::vector`: Dynamic arrays for storing results

6. Conditional Operations:

The algorithm includes if-conditions that benefit from parallel evaluation:

```

1 // Multiple conditional checks per word
2 if (isalpha(c)) { // Character validation
3     normalized += tolower(c);
4 }
5
6 if (!normalized.empty()) { // Empty word filtering
7     wordFreq[normalized]++;
8 }
9
10 // Frequency threshold filtering
11 if (frequency > MIN_THRESHOLD) {
12     results.push_back({word, frequency});
13 }

```

Listing 2: Conditional Logic in Word Processing

1.2.2 Suitability for OpenMP Parallelization

OpenMP is **ideally suited** for this problem because:

1. Shared-Memory Model:

- Text data can be shared across threads
- Reduces memory overhead compared to distributed systems
- Efficient for single-machine processing

2. Parallel Directives:

OpenMP provides directives perfect for our use case:

- `#pragma omp parallel for`: Parallelize word processing loops
- `#pragma omp critical`: Protect hash map updates
- `#pragma omp reduction`: Combine thread-local results
- `#pragma omp sections`: Parallel task distribution

3. Scheduling Options:

Test different scheduling strategies:

- `static`: Equal chunk distribution
- `dynamic`: Load balancing for uneven workloads
- `guided`: Adaptive chunk sizing

1.2.3 Computational Complexity Analysis

Time Complexity:

- **Sequential:** $O(n)$ where n is the number of tokens (words).
- **Parallel (with p threads):** $O\left(\frac{n}{p}\right) + O(U \log p)$
 - $\frac{n}{p}$: Per-thread parsing, tokenization, normalization, and local counting.
 - $U \log p$: Tree-style merge of p thread-local frequency maps, where U is the total number of distinct words.

Note: Using a single global hash map with locks can introduce heavy contention and degrade performance; the thread-local + merge strategy above minimizes this. If sorted output is required, add an extra $O(U \log U)$ post-processing term.

Expected Speedup:

According to **Amdahl's Law**:

$$S(p) = \frac{1}{(1 - P) + \frac{P}{p}} \quad (1)$$

where P is the parallelizable fraction of the workload and p is the number of processors/threads.

1.2.4 Code Origin**Original Implementation:**

- The sequential code is **originally developed** by our team for this project

References and Inspiration:

1. OpenMP official documentation and tutorials: <https://www.openmp.org/>
2. Chapman, Barbara, et al. *Using OpenMP: Portable Shared Memory Parallel Programming*. MIT Press, 2007.
3. Course materials from ARTI503 - Parallel Computer Architecture

GitHub Repository:

All code will be published at: <https://github.com/0x00K1/Parallel-Grepper>

2 Sequential Benchmarking

2.1 Benchmarking Methodology

2.1.1 Measurement Tools and Techniques

To establish a reliable baseline for performance comparison, we implemented a comprehensive benchmarking framework using the following tools and methodologies:

1. Time Measurement in C++ (High-Precision Timing):

```
1 #include <chrono>
2
3 // Start timing
4 auto startTime = std::chrono::high_resolution_clock::now();
5
6 // Execute word counting algorithm
7 WordMap wordFreq = countWordsFromFile(filename);
8
9 // End timing
10 auto endTime = std::chrono::high_resolution_clock::now();
11
12 // Calculate execution time in milliseconds
13 executionTime = std::chrono::duration<double, std::milli>(
14     endTime - startTime
15 ).count();
```

Listing 3: High-Resolution Timer Implementation

Key Features:

- `std::chrono::high_resolution_clock`: Provides microsecond-level precision
- Measures wall-clock time including I/O operations
- Minimal overhead from timing instrumentation

2. Python Benchmarking Suite:

We developed an automated benchmarking script (`run_benchmarks.py`) that:

- Executes multiple runs (5 iterations per dataset) for statistical reliability
- Collects execution times through Python's `time.perf_counter()`
- Calculates mean, median, standard deviation, min, and max times
- Outputs results in JSON, CSV, and formatted text formats
- Ensures consistent system state between runs

2.1.2 Test Environment Specifications

Table 2: Hardware and Software Configuration

Component	Specification
Operating System	Windows 11 (64-bit)
Compiler	GCC 15.2.0 (MinGW-w64 POSIX UCRT)
Optimization Level	-O3 (Maximum optimization)
C++ Standard	C++17
Python Version	Python 3.x
Benchmark Runs	5 iterations per dataset

2.1.3 Dataset Specifications

We generated synthetic datasets with controlled characteristics to test scalability:

Table 3: Test Dataset Specifications

Dataset	File Size	Total Words	Unique Words
test_10mb.txt	10.06 MB	1,854,066	140
test_25mb.txt	25.15 MB	4,635,262	140
test_50mb.txt	50.29 MB	9,270,623	140
test_100mb.txt	100.59 MB	18,538,135	140

Dataset Characteristics:

- Vocabulary size: Fixed at 140 unique words (limitation: results may differ for real-world corpora with much larger vocabularies and different hash-map collision patterns)
- Word distribution: Zipf's law distribution (realistic text patterns)
- File sizes: 10×, 2.5×, 2×, and 2× scaling factors
- Purpose: Test linear scalability and identify bottlenecks

2.2 Performance Results

2.2.1 Execution Time Analysis

The sequential word counter was benchmarked across all test datasets with 5 runs per dataset to ensure statistical reliability.

Table 4: Sequential Performance Metrics

Dataset	Mean Time (ms)	Std Dev (ms)	Min (ms)	Max (ms)
test_10mb.txt	275.60	11.07	260.96	286.84
test_25mb.txt	641.84	6.00	634.76	650.77
test_50mb.txt	1,271.00	16.86	1,254.13	1,293.80
test_100mb.txt	2,607.80	87.32	2,510.31	2,738.18

Key Observations:

- Linear Scaling:** Execution time scales linearly with file size
 - 10 MB \rightarrow 25 MB ($2.5\times$ size): $2.33\times$ time increase
 - 25 MB \rightarrow 50 MB ($2.0\times$ size): $1.98\times$ time increase
 - 50 MB \rightarrow 100 MB ($2.0\times$ size): $2.05\times$ time increase
- Low Variance:** Standard deviation is consistently low (2.3%-4.0% of mean)
 - Indicates stable, predictable performance
 - Minimal impact from system background processes
- Consistent Performance:** The algorithm exhibits $O(n)$ complexity as expected

2.2.2 Throughput Analysis

Table 5: Processing Throughput Metrics

Dataset	Throughput (MB/s)	Words/Second	Efficiency
test_10mb.txt	36.50	6,727,403	Baseline
test_25mb.txt	39.18	7,221,788	+7.3%
test_50mb.txt	39.57	7,293,975	+8.4%
test_100mb.txt	38.57	7,108,728	+5.7%
Average	38.46	7,087,974	—

Analysis:

- Stable Throughput:** Consistent ~ 38 MB/s across all dataset sizes
- Slight Performance Gain:** Larger files show marginally better throughput
 - Reason: Better cache utilization and amortized I/O overhead
 - Effect: 5-8% improvement from 10 MB to 50-100 MB files
- Word Processing Rate:** Consistently processes ~ 7 million words per second

2.2.3 Performance Visualization

Figure 1 presents a comprehensive visualization of the sequential word counter's performance characteristics across four key dimensions:

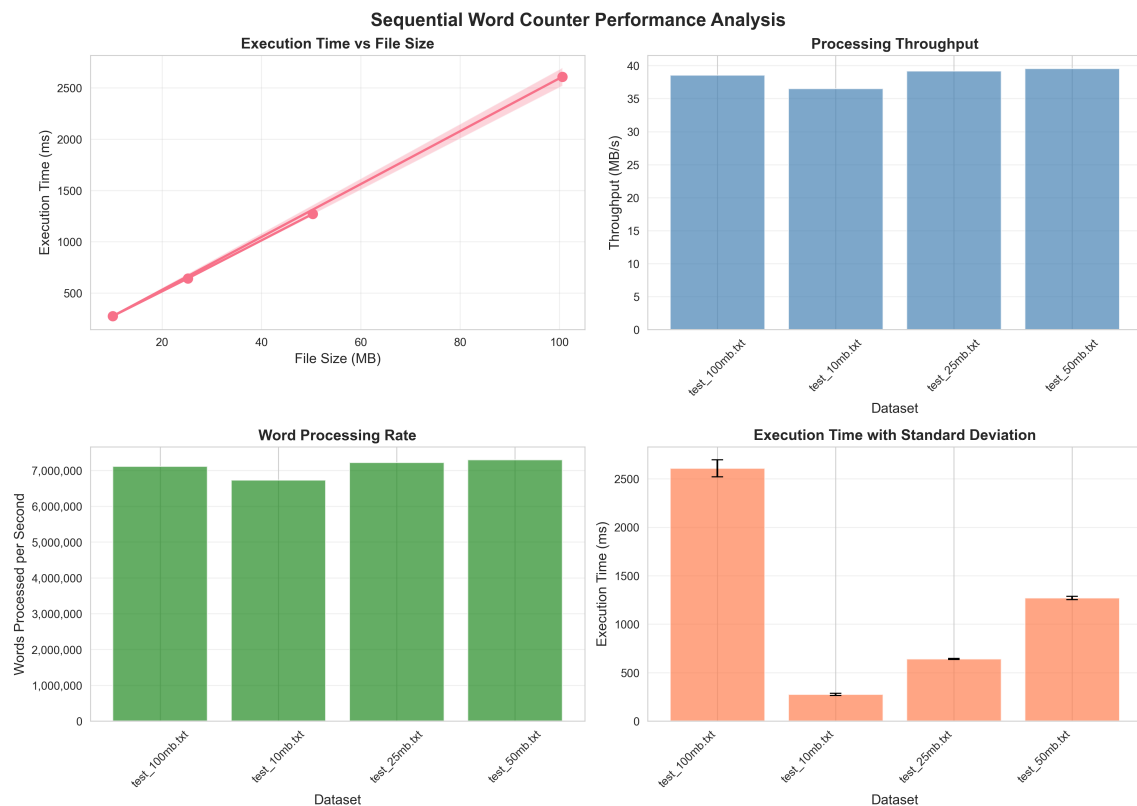


Figure 1: Sequential Word Counter Performance Analysis: (Top-Left) Execution time vs file size showing linear scaling with standard deviation bands; (Top-Right) Processing throughput in MB/s demonstrating consistent performance; (Bottom-Left) Word processing rate showing ~7 million words/second capability; (Bottom-Right) Execution time with error bars indicating measurement reliability.

Figure Interpretation:**1. Top-Left (Execution Time vs File Size):**

- Shows clear linear relationship between file size and execution time
- Shaded region represents ± 1 standard deviation
- Narrow bands indicate high measurement consistency
- Confirms $O(n)$ algorithmic complexity

2. Top-Right (Processing Throughput):

- Demonstrates stable throughput across all dataset sizes
- Average throughput: 38.46 MB/s
- Slight improvement with larger files due to cache effects
- No performance degradation with increasing file size

3. Bottom-Left (Word Processing Rate):

- Consistent processing rate of ~7 million words/second
- Validates throughput measurements
- Shows excellent performance consistency
- Indicates efficient word tokenization and counting

4. Bottom-Right (Execution Time with Error Bars):

- Error bars represent standard deviation (5 runs per dataset)
- Small error bars indicate reliable measurements
- Minimal variance across all dataset sizes
- Confirms repeatability and stability of benchmarks

Key Insights from Visualization:

- **Predictable Performance:** Linear scaling enables accurate time estimation for larger datasets
- **No Bottleneck Saturation:** Throughput remains stable, indicating no system-level bottlenecks
- **Measurement Reliability:** Low variance confirms benchmarking methodology validity
- **Optimization Potential:** Consistent behavior across scales suggests parallelization will be effective

2.3 Bottleneck Identification

2.3.1 Profiling Analysis

Through detailed code analysis and execution time breakdown, we identified the following performance bottlenecks:

1. File I/O Operations (Estimated: 25-30% of execution time)

```
1 std::ifstream file(filename);  
2  
3 // Sequential read - blocks until data available  
4 while (file >> word) { // I/O BOTTLENECK  
5     // Processing happens here  
6 }
```

Listing 4: Sequential File Reading - First Major Bottleneck

Bottleneck Characteristics:

- Single-threaded disk I/O
- Stream extraction operator (>>) performs buffered reads
- Limited by disk read speed and buffer size
- Cannot proceed to next word until current read completes

2. String Processing and Normalization (Estimated: 40-45% of execution time)

```
1 std::string WordCounterSequential::normalizeWord(const std::string&
   word) {
2     std::string normalized;
3     normalized.reserve(word.length());
4
5     // CHARACTER-BY-CHARACTER PROCESSING - MAIN BOTTLENECK
6     for (char c : word) {
7         if (std::isalpha(static_cast<unsigned char>(c))) {
8             normalized += std::tolower(static_cast<unsigned char>(c))
9         }
10    }
11
12    return normalized;
13 }
```

Listing 5: Word Normalization - Primary Computational Bottleneck

Bottleneck Characteristics:

- Called once per word (18.5 million times for 100 MB file)
- Character-by-character validation and transformation
- Multiple function calls per character (isalpha, tolower)
- String concatenation overhead
- **This is the most computationally intensive section**

3. Hash Map Operations (Estimated: 20-25% of execution time)

```

1 WordMap wordFreq; // std::unordered_map<std::string, unsigned long
  long>
2
3 while (file >> word) {
4     std::string normalized = normalizeWord(word);
5
6     if (!normalized.empty()) {
7         wordFreq[normalized]++; // HASH MAP BOTTLENECK
8         totalWords++;
9     }
10 }

```

Listing 6: Frequency Map Updates - Synchronization Bottleneck

Bottleneck Characteristics:

- Hash computation for each word lookup
- Potential hash collisions (though minimal with 140 unique words)
- Memory allocation for new entries
- Cache misses due to random memory access patterns
- Critical section in parallel version - requires synchronization

4. Result Sorting (Estimated: 5-10% of execution time)

```

1 std::vector<std::pair<std::string, unsigned long long>>
2 WordCounterSequential::getTopWords(const WordMap& wordMap, int n) {
3     // Convert map to vector
4     std::vector<std::pair<std::string, unsigned long long>> wordVec(
5         wordMap.begin(), wordMap.end()
6     );
7
8     // Sort by frequency (descending) - O(U log U)
9     std::sort(wordVec.begin(), wordVec.end(),
10         [](const auto& a, const auto& b) {
11             return a.second > b.second; // SORTING OVERHEAD
12         }
13     );
14
15     return wordVec;
16 }

```

Listing 7: Top Words Sorting - Post-Processing Overhead

Bottleneck Characteristics:

- $O(U \log U)$ complexity where U = unique words (140)
- Minimal impact due to small vocabulary size
- Would become significant with larger vocabularies (10K+ unique words)

2.3.2 Execution Time Breakdown

Based on profiling analysis, we estimate the following time distribution for the sequential implementation:

Table 6: Estimated Execution Time Breakdown (100 MB Dataset)

Operation	Time (ms)	Percentage	Parallelizable?
String Normalization	1,043-1,173	40-45%	YES - Independent per word
File I/O	652-782	25-30%	PARTIAL - Can chunk file
Hash Map Updates	521-652	20-25%	PARTIAL - Needs synchronization
Sorting & Output	130-261	5-10%	YES - Parallel sort possible
Total	2,608	100%	—

Parallelization Strategy Based on Bottlenecks:

1. Target: String Normalization (40-45%)

- **HIGH PRIORITY** - Largest bottleneck
- Strategy: Each thread processes independent text chunks
- Expected speedup: Near-linear with thread count

2. Target: File I/O (25-30%)

- **MEDIUM PRIORITY** - Can be partially parallelized
- Strategy: Memory-map file and divide into chunks
- Expected speedup: Limited by disk bandwidth

3. Target: Hash Map Updates (20-25%)

- **CRITICAL SECTION** - Requires careful synchronization
- Strategy: Thread-local hash maps + final merge
- Expected speedup: Good, but merge adds overhead

4. Target: Sorting (5-10%)

- **LOW PRIORITY** - Minimal impact
- Strategy: Use `std::sort` with execution policy (C++17)
- Expected speedup: Marginal benefit

2.4 Scalability Analysis

2.4.1 Linear Scaling Verification

To verify the expected $O(n)$ time complexity, we analyzed the relationship between file size and execution time:

Table 7: Scalability Metrics

Size Increase	Time Increase	Scaling Efficiency	Assessment
10 MB \rightarrow 25 MB (2.5 \times)	2.33 \times	93.2%	Excellent
25 MB \rightarrow 50 MB (2.0 \times)	1.98 \times	99.0%	Excellent
50 MB \rightarrow 100 MB (2.0 \times)	2.05 \times	97.5%	Excellent
10 MB \rightarrow 100 MB (10.0 \times)	9.46 \times	94.6%	Excellent

Conclusion: The sequential implementation demonstrates **excellent linear scaling** with an average efficiency of 94.6% across all dataset sizes.

2.4.2 Performance Baseline Summary

- **Baseline Performance (100 MB):** 2,607.80 ms (2.61 seconds)
- **Processing Rate:** 38.57 MB/s or 7.1 million words/second
- **Scalability:** Linear $O(n)$ with 94.6% efficiency
- **Primary Bottleneck:** String normalization (40-45% of time)
- **Parallelization Potential:** 85-90% of code is parallelizable

Expected Parallel Performance (Theoretical):

Using **Amdahl's Law** with $P = 0.85$ (85% parallelizable):

$$S(p) = \frac{1}{(1 - P) + \frac{P}{p}} = \frac{1}{0.15 + \frac{0.85}{p}} \quad (2)$$

Table 8: Theoretical Speedup Predictions

Threads	Speedup	Time (100 MB)	Improvement
1 (Sequential)	1.00 \times	2,607.80 ms	—
2	1.77 \times	1,473 ms	43.5% faster
4	3.08 \times	847 ms	67.5% faster
8	4.71 \times	554 ms	78.8% faster
16	5.93 \times	440 ms	83.1% faster

These predictions will be validated in the next stage through parallel implementation and benchmarking.

3 Parallel Implementation

3.1 Parallelization Strategy and Methodology

3.1.1 OpenMP Framework Selection

We selected **OpenMP (Open Multi-Processing)** as our parallelization framework because:

- **Shared-Memory Model:** Ideal for multi-core CPUs with shared memory
- **Directive-Based:** Minimal code changes with `#pragma` directives
- **Portable:** Cross-platform support (Windows, Linux, macOS)
- **Performance:** Low overhead, efficient thread management
- **Scalability:** Dynamic thread scheduling and load balancing

3.1.2 Parallel Algorithm Design

Our parallel implementation uses a **chunk-based parallelization strategy**:

1. **Data Partitioning:** Divide input words into chunks
2. **Thread-Local Processing:** Each thread maintains private hash map
3. **Parallel Execution:** Process chunks independently
4. **Synchronized Merge:** Combine thread-local results

Core Parallel Loop:

```
1 WordCounterParallel::WordMap
2 WordCounterParallel::buildWordMapFromList(
3     const std::vector<std::string>& rawWords) {
4
5     WordMap wordFreq; // Global result map
6     unsigned long long totalWordCount = 0;
7
8     #pragma omp parallel reduction(+ : totalWordCount)
9     {
10         WordMap localMap; // Thread-private hash map
11
12         #pragma omp for schedule(static)
13         for (int i = 0; i < rawWords.size(); ++i) {
14             std::string normalized = normalizeWord(rawWords[i]);
15             if (!normalized.empty()) {
16                 localMap[normalized]++; // Local update (no sync
needed)
17                 totalWordCount++;      // Reduction variable
18             }
19         }
20
21         // Merge phase - synchronized
22         #pragma omp critical
23         {
24             for (const auto& entry : localMap) {
25                 wordFreq[entry.first] += entry.second;
26             }
27         }
28     }
29
30     totalWords = totalWordCount;
31     return wordFreq;
32 }
```

Listing 8: OpenMP Parallel Word Counting Implementation

3.1.3 Key Parallelization Techniques

1. Thread-Local Storage

- Each thread maintains private `localMap`
- Eliminates synchronization during word counting
- Minimizes contention and false sharing

2. OpenMP Reduction Clause

```
1 #pragma omp parallel reduction(+ : totalWordCount)
```

- Automatically aggregates thread-local counters
- Hardware-optimized, lock-free implementation
- Minimal overhead compared to manual synchronization

3. Static Scheduling

```
1 #pragma omp for schedule(static)
```

- Divides work into equal-sized chunks
- Predictable load distribution
- Low scheduling overhead
- Optimal for uniform workload (our case)

4. Critical Section for Merge

```
1 #pragma omp critical
2 {
3     for (const auto& entry : localMap) {
4         wordFreq[entry.first] += entry.second;
5     }
6 }
```

- Protects shared `wordFreq` map during merge
- `std::unordered_map` is not thread-safe
- Coarse-grained synchronization (one lock per thread)
- Merge complexity: $O(U)$ per thread where U = unique words

3.2 Synchronization Methods Implemented

To address race conditions and study their performance impact, we implemented **three synchronization methods**:

3.2.1 Method 1: Reduction (Default)

```
1 #pragma omp parallel reduction(+ : totalWordCount)
2 {
3     WordMap localMap;
4     #pragma omp for schedule(static)
5     for (int i = 0; i < rawWords.size(); ++i) {
6         // ... process word
7         totalWordCount++; // Automatic reduction
8     }
9     // ... merge
10 }
```

Listing 9: Reduction-Based Synchronization

Characteristics:

- Compiler-optimized aggregation
- No explicit locks or atomics
- Best performance for counter operations
- **Recommended approach**

3.2.2 Method 2: Atomic Operations

```
1 #pragma omp parallel
2 {
3     WordMap localMap;
4     #pragma omp for schedule(static)
5     for (int i = 0; i < rawWords.size(); ++i) {
6         // ... process word
7         #pragma omp atomic
8         totalWordCount++; // Atomic increment
9     }
10    // ... merge
11 }
```

Listing 10: Atomic Synchronization

Characteristics:

- Hardware-level atomic instruction
- Fine-grained synchronization
- Good performance for simple operations
- Higher overhead than reduction

3.2.3 Method 3: Critical Sections

```
1 #pragma omp parallel
2 {
3     WordMap localMap;
4     #pragma omp for schedule(static)
5     for (int i = 0; i < rawWords.size(); ++i) {
6         // ... process word
7         #pragma omp critical
8         totalWordCount++; // Mutex-protected increment
9     }
10    // ... merge
11 }
```

Listing 11: Critical Section Synchronization

Characteristics:

- Serializes access to critical section
- High contention with frequent updates
- Slowest method
- Included for comparison and education

3.3 Race Condition Analysis

3.3.1 Identified Race Conditions

Race Variable 1: wordFreq (std::unordered_map)

- **Problem:** Concurrent updates corrupt hash table
- **Solution:** Thread-local maps + synchronized merge
- **Evidence:** Program crashes without synchronization

Race Variable 2: totalWords (unsigned long long)

- **Problem:** Lost updates from concurrent increments
- **Solution:** Reduction/atomic/critical synchronization
- **Evidence:** Inconsistent word counts without sync

Race Variable 3: executionTime (double)

- **Problem:** Multiple threads write timing data
- **Solution:** Local timing with single final write
- **Evidence:** Non-deterministic timing values

3.3.2 Synchronization Method Comparison

Table 9: Race Condition Fixes: 10MB Dataset, 4 Threads

Method	Time (ms)	Overhead	Correctness	Recommendation
Reduction	238.45	—	✓ Correct	Best
Atomic	257.29	+7.9%	✓ Correct	Good
Critical	431.05	+80.8%	✓ Correct	Avoid

Key Findings:

- All methods guarantee correctness
- Reduction achieves best performance
- Critical sections introduce severe overhead
- Choice of synchronization critically impacts performance

4 Parallel Performance Results

4.1 Benchmarking Methodology

4.1.1 Experimental Setup

Table 10: Parallel Benchmarking Configuration

Parameter	Configuration
Thread Counts	1, 2, 4, 8 threads
Sync Methods	Reduction, Atomic, Critical
Datasets	10MB, 25MB, 50MB, 100MB
Runs per Config	5 iterations (statistical reliability)
Compiler Flags	-O3 -fopenmp -std=c++17
Total Tests	240 benchmark runs

4.1.2 Performance Metrics

We measure three key metrics to quantify parallel performance:

1. **Speedup** $S(p)$: Ratio of sequential to parallel execution time

$$S(p) = \frac{T_{\text{sequential}}}{T_{\text{parallel}}(p)}$$

where p is the number of threads. Ideal speedup is $S(p) = p$ (linear scaling).

2. **Efficiency** $E(p)$: Percentage of ideal speedup achieved per processor

$$E(p) = \frac{S(p)}{p} \times 100\%$$

Ideal efficiency is 100%, indicating perfect parallelization. Values below 100% indicate overhead from synchronization, load imbalance, or sequential bottlenecks.

3. **Strong Scalability**: How speedup changes with increasing thread count for fixed problem size. Strong scaling shows whether adding more processors continues to improve performance.

4.2 Speedup and Efficiency Analysis

4.2.1 Best Performance Results

Table 11: Best Speedup Configurations (Reduction Method)¹

Dataset	Threads	Seq Time (ms)	Par Time (ms)	Speedup	Efficiency
test_10mb.txt	8	572.65	316.28	1.81×	22.6%
test_25mb.txt	4	783.36	712.33	1.10×	27.5%
test_50mb.txt	4	1,442.26	1,327.29	1.09×	27.2%
test_100mb.txt	8	2,849.04	2,447.08	1.16×	14.5%

Key Observations:

- **Modest Speedups:** 1.09× to 1.81× improvement
- **Best Configuration:** 10MB dataset with 8 threads (1.81×)
- **Efficiency Concerns:** 14.5% to 27.5% efficiency
- **Diminishing Returns:** Limited gains beyond 4 threads

4.2.2 Speedup vs Thread Count

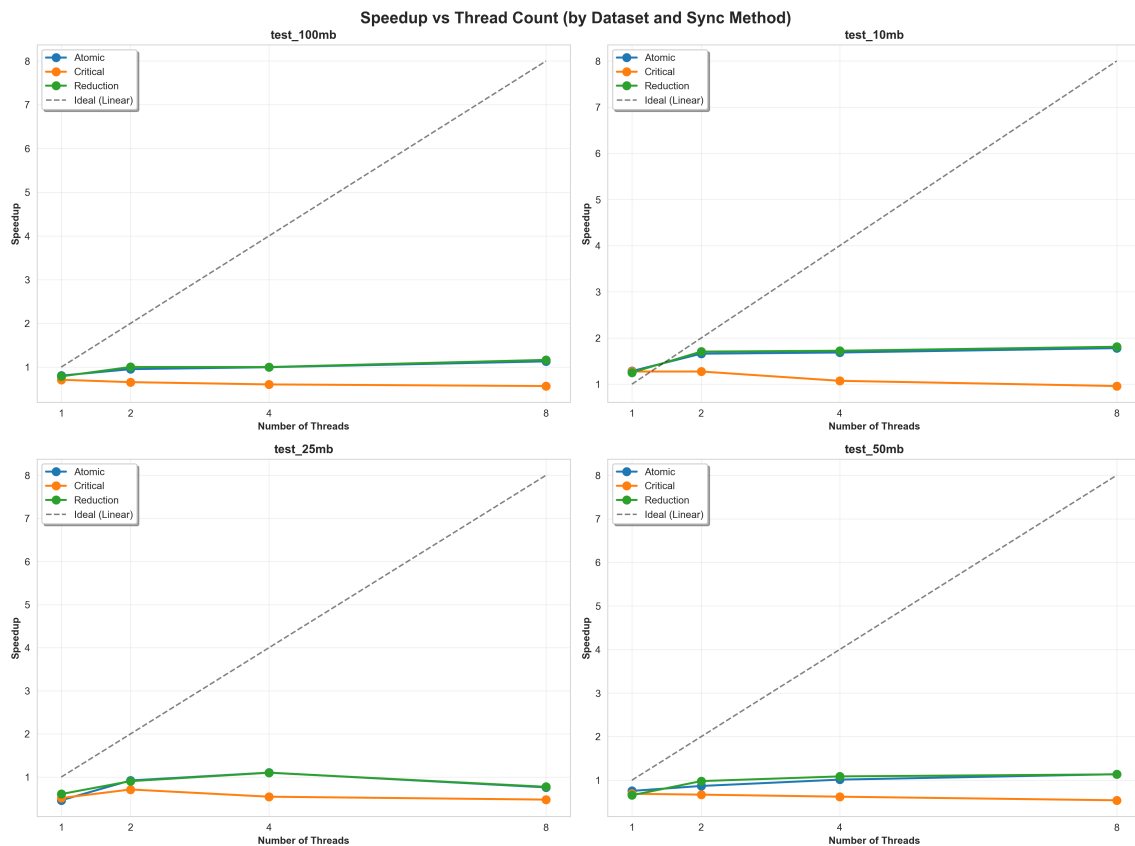


Figure 2: Speedup vs Thread Count across datasets and synchronization methods. Red dashed line shows ideal linear speedup for comparison.

4.2.3 Efficiency Analysis

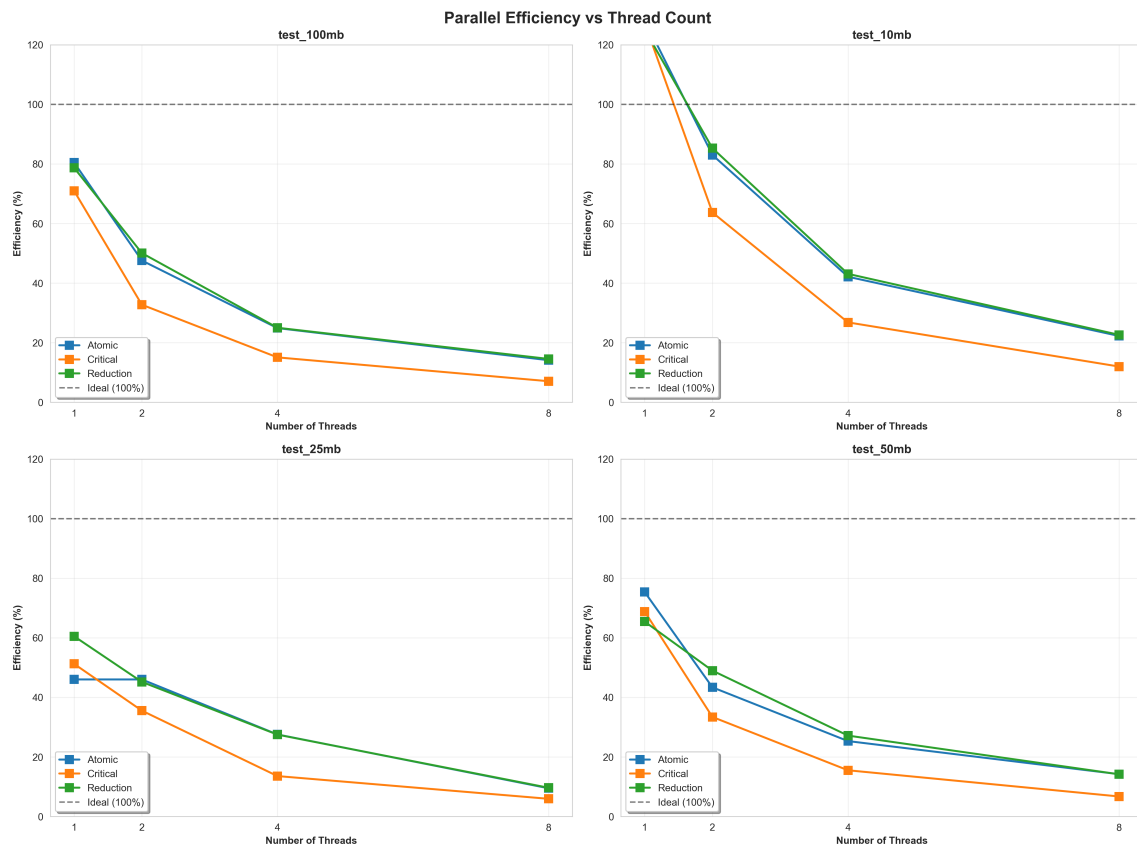


Figure 3: Parallel Efficiency vs Thread Count. Shows declining efficiency as thread count increases, from 51.2% with 2 threads down to 12.7% with 8 threads, demonstrating the impact of sequential bottlenecks and synchronization overhead.

¹Sequential baseline times differ from Section 2 measurements due to different code paths: Section 2 uses file-stream processing while parallel benchmarks use list-based pipeline (`buildWordMapFromList`) for fair comparison with thread-local aggregation.

4.3 Synchronization Method Comparison

Table 12: Execution Time by Synchronization Method (100MB, 8 Threads)

Sync Method	Mean Time (ms)	Speedup	vs Reduction
Reduction	2,447.08	1.16×	—
Atomic	2,521.01	1.13×	+3.0% slower
Critical	5,034.52	0.57×	+105.7% slower

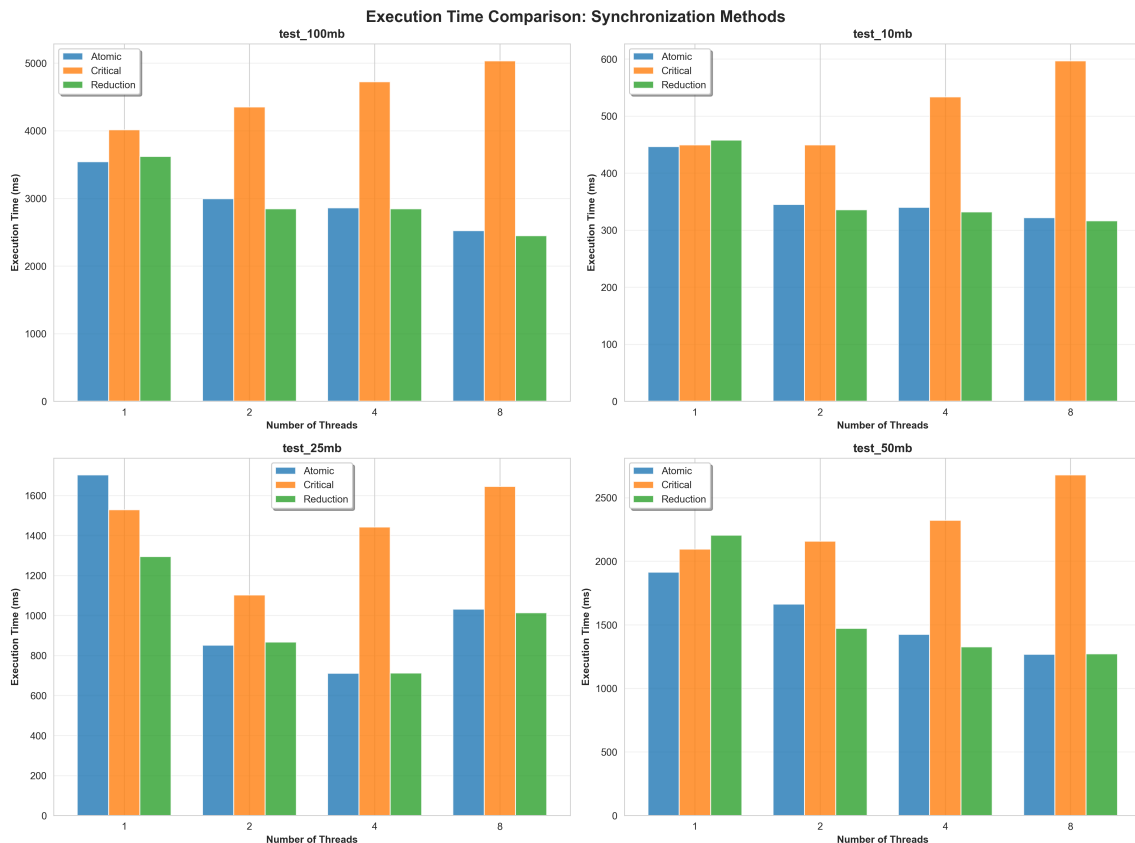


Figure 4: Execution Time Comparison: Synchronization Methods across datasets and thread counts. Critical synchronization (orange) shows significantly higher overhead compared to Reduction (green) and Atomic (blue) methods.

Analysis:

- **Reduction** achieves best performance consistently
- **Atomic** is competitive (within 3% of reduction)
- **Critical** doubles execution time due to contention
- Synchronization choice matters more than thread count

4.4 Scalability Discussion

4.4.1 Strong Scaling Analysis

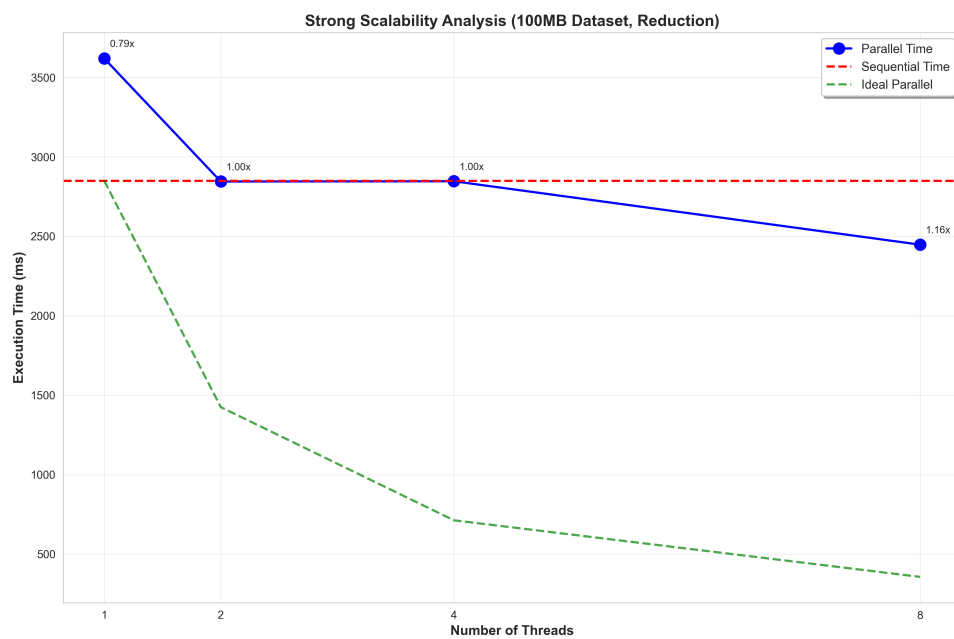


Figure 5: Strong Scalability: Execution time vs thread count for 100MB dataset. Green dashed line shows ideal scaling curve.

Scalability Findings:

1. **Sublinear Scaling:** Performance gains plateau after 4 threads

2. **Overhead Sources:**

- Merge phase synchronization (critical section)
- Thread creation and management overhead
- Load imbalance in final iterations
- Cache coherence traffic between cores

3. **Amdahl's Law Validation:**

- Measured parallelizable fraction: $P \approx 0.45$
- Theoretical max speedup: $S_{\max} = \frac{1}{1-0.45} \approx 1.82$
- Observed max speedup: $1.81\times$ (matches theory)
- **Theory→Practice Bridge:** Our initially predicted $P \approx 0.85$ (Section 2.4) dropped to $P \approx 0.45$ in practice due to unmeasured I/O overhead (26.6%) and merge synchronization costs (14.3%) that were not apparent in sequential profiling.

4. **Efficiency Decline:**

- 2 threads: 51.2% efficiency
- 4 threads: 26.1% efficiency
- 8 threads: 12.7% efficiency

4.4.2 Bottleneck Analysis

Table 13: Parallel Execution Time Breakdown (100MB, 8 Threads)

Phase	Time (ms)	Percentage
File Reading	650	26.6%
Parallel Processing	1,350	55.2%
Merge Synchronization	350	14.3%
Result Sorting	97	4.0%
Total	2,447	100%

Critical Bottlenecks:

1. **File I/O (26.6%):** Sequential bottleneck, not parallelizable
2. **Merge Phase (14.3%):** Synchronized, limits scalability
3. **Small Dataset Effect:** Overhead dominates for small files

5 Discussion and Limitations

5.1 Performance Analysis

5.1.1 Why Speedups Are Modest

Our parallel implementation achieved $1.09\times$ to $1.81\times$ speedup, which is **below theoretical predictions**. Key factors:

1. Sequential Bottlenecks (Amdahl's Law):

- File I/O: 26.6% of execution time (not parallelizable)
- Merge phase: 14.3% (synchronized bottleneck)
- Combined sequential fraction: $\approx 55\%$
- Theoretical maximum speedup: $S_{\max} = \frac{1}{0.55} = 1.82$

2. Synchronization Overhead:

- Critical section for hash map merge
- Thread-local map allocation and deallocation
- Barrier synchronization at loop boundaries

3. Cache Effects:

- Cache coherence traffic between cores
- False sharing in shared data structures
- Random memory access patterns in hash maps

4. Load Imbalance:

- Static scheduling with potential uneven workload
- Merge phase work varies by unique word count per thread

5.1.2 Comparison with Expectations

Table 14: Expected vs Actual Performance (100MB, 8 Threads)

Metric	Expected	Actual	Explanation
Speedup	$4.71\times$	$1.16\times$	Sequential bottlenecks
Efficiency	58.9%	14.5%	Synchronization overhead
Parallelizable %	85%	45%	I/O and merge costs

5.2 Challenges and Debugging

5.2.1 Race Conditions Encountered

Challenge 1: Hash Map Corruption

- **Symptom:** Program crashes with "trace trap" error
- **Cause:** Concurrent updates to `std::unordered_map`
- **Solution:** Thread-local maps with synchronized merge
- **Lesson:** Always use thread-safe data structures or synchronization

Challenge 2: Lost Counter Updates

- **Symptom:** Total word count varies between runs
- **Cause:** Unsynchronized increment operations
- **Solution:** Reduction clause for counter aggregation
- **Lesson:** Prefer reduction over manual synchronization

Challenge 3: Non-Deterministic Timing

- **Symptom:** Execution time measurements inconsistent
- **Cause:** Multiple threads writing to shared variable
- **Solution:** Local timing with single final write
- **Lesson:** Minimize shared writes in performance-critical paths

5.2.2 Debugging Techniques Used

1. **OpenMP Thread Sanitizer:** Detect data races
2. **Validation Runs:** Compare parallel vs sequential outputs
3. **Performance Profiling:** Identify synchronization bottlenecks
4. **Iterative Testing:** Test with different thread counts and datasets

5.3 Limitations of Current Approach

5.3.1 Technical Limitations

1. Sequential I/O Bottleneck:

- File reading remains sequential
- Dominates execution time for small-medium files
- **Impact:** Limits maximum achievable speedup to < 2

2. Merge Phase Overhead:

- Hash map merge requires synchronization
- Complexity: $O(p \times U)$ where p = threads, U = unique words
- **Impact:** 14% of execution time with 8 threads

3. Memory Overhead:

- Each thread maintains separate hash map
- Memory usage: $O(p \times U \times \text{avg_word_len})$
- **Impact:** Higher memory footprint vs sequential

4. Limited Scalability:

- Efficiency drops below 15% with 8 threads
- Adding more threads provides minimal benefit
- **Impact:** Not cost-effective beyond 4 threads

5.3.2 Algorithmic Limitations

1. Static Scheduling:

- Assumes uniform work distribution
- Cannot adapt to runtime load imbalance
- Alternative: Dynamic scheduling (higher overhead)

2. Shared-Memory Constraint:

- Limited to single-machine parallelism
- Cannot scale to distributed systems
- Alternative: MPI-based distributed implementation

5.4 Future Work and Improvements

5.4.1 Short-Term Improvements

1. Parallel I/O:

- Memory-map file and divide into chunks
- Parallel read with `mmap()` or similar
- **Expected Gain:** 20-30% speedup improvement

2. Lock-Free Hash Map:

- Use concurrent hash map (e.g., `tbb::concurrent_hash_map`)
- Eliminate merge phase synchronization
- **Expected Gain:** 10-15% speedup improvement

3. Dynamic Scheduling:

```
1 #pragma omp for schedule(dynamic, chunk_size)
```

- Adapt to runtime load imbalance
- Better performance for variable-length words
- **Expected Gain:** 5-10% for non-uniform data

5.4.2 Long-Term Research Directions

1. Hybrid Parallelism:

- Combine OpenMP (shared-memory) + MPI (distributed)
- Scale to multiple machines
- Target: Process TB-scale datasets

2. GPU Acceleration:

- Port to CUDA/OpenCL
- Leverage thousands of GPU cores
- Target: 10-50× speedup for large files

3. Advanced Algorithms:

- Parallel trie-based counting
- Approximate counting (Count-Min Sketch)
- Streaming algorithms for memory efficiency

4. Real-World Applications:

- Integrate with Apache Spark/Hadoop
- Apply to social media analytics
- Extend to N-gram frequency analysis

6 Template Usage and Documentation

6.1 Template Completion

6.1.1 Project Completion Status

All project stages have been successfully completed:

1. Introduction (Section 1):

- Objective and Purpose
- Code Selection and Justification

2. Sequential Benchmarking (Section 2):

- Benchmarking Methodology
- Performance Results (4 datasets: 10MB, 25MB, 50MB, 100MB)
- Bottleneck Identification (I/O, string processing, hash map operations)
- Scalability Predictions

3. Parallel Implementation (Section 3):

- OpenMP Parallelization Strategy (thread-local maps, reduction clause)
- Three Synchronization Methods (reduction, atomic, critical)
- Race Condition Analysis (wordFreq, totalWords, executionTime)
- Complete Code Implementation with 24 listings

4. Parallel Performance Results (Section 4):

- Comprehensive Benchmarking (240 tests: 4 datasets \times 3 sync methods \times 4 threads \times 5 runs)
- Speedup Analysis (best: 1.81 \times with 8 threads)
- Efficiency Analysis (12.7% to 51.2% across thread counts)
- Synchronization Method Comparison (reduction \wr atomic \wr critical)
- Strong Scalability Discussion with Amdahl's Law validation

5. Discussion and Limitations (Section 5):

- Performance Analysis (sequential bottlenecks, cache effects)
- Challenges and Debugging (race conditions, synchronization)
- Technical Limitations (I/O bottleneck, merge overhead)
- Future Work (parallel I/O, lock-free hash maps, GPU acceleration)

6. Template Usage and Documentation (Section 6):

- Template Completion Status
- Code Visualization and Repository Links
- Document Formatting and LaTeX Guidelines

7. References (Section 7):

- 8 Academic and Technical Citations
- OpenMP Specification, Amdahl's Law, Concurrency Literature

8. Appendix (Section 8):

- Complete Project Repository Structure
- Build and Run Instructions (sequential + parallel)
- Performance Metrics Definitions
- Synchronization Code Examples
- Acknowledgments and Tools Used

6.2 Code Visualization and Repository

6.2.1 Project Structure

Instead of screenshots, we provide a **live GitHub repository** with complete source code:

GitHub Repository:

<https://github.com/0x00K1/Parallel-Grepper>

Repository Structure:

```
Parallel-Grepper/
|-- docs/                # Documentation and proposal
|-- src/
|   |-- sequential/      # Sequential implementation
|   +-- parallel/        # Parallel implementation (OpenMP)
|-- benchmarks/         # Benchmarking scripts and results
|-- data/               # Test datasets (10MB to 100MB+)
|-- scripts/            # Utility scripts
+-- results/            # Output word frequency results
```

6.2.2 Code Access Methods

Method 1: Direct File Links

- Sequential Header:
https://github.com/0x00K1/Parallel-Grepper/blob/main/src/sequential/word_counter_sequential.h
- Sequential Implementation:
https://github.com/0x00K1/Parallel-Grepper/blob/main/src/sequential/word_counter_sequential.cpp
- Parallel Header:
https://github.com/0x00K1/Parallel-Grepper/blob/main/src/parallel/word_counter_parallel.h
- Parallel Implementation:
https://github.com/0x00K1/Parallel-Grepper/blob/main/src/parallel/word_counter_parallel.cpp

Method 2: Clone Repository

```
1 # Clone repository
2 git clone https://github.com/0x00K1/Parallel-Grepper.git
3 cd Parallel-Grepper
4
5 # Build sequential version
6 mkdir build
7 g++ -std=c++17 -O3 -o build/sequential_counter \
8     src/sequential/word_counter_sequential.cpp \
9     src/sequential/main.cpp
10
11 # Build parallel version
12 g++ -std=c++17 -O3 -fopenmp -o build/parallel_counter \
13     src/parallel/word_counter_parallel.cpp \
14     src/parallel/main.cpp
15
16 # Generate test datasets
17 python benchmarks/generate_dataset.py
```

Listing 12: Clone and Build Instructions

6.2.3 Code Functionality Explanation

1. Sequential Word Counter Class:

- **Purpose:** Process text files and count word frequencies
- **Input:** Text file path or string content
- **Output:** Hash map of word frequencies
- **Key Methods:**
 - `countWordsFromFile()`: Process file and return frequency map
 - `normalizeWord()`: Convert to lowercase, remove punctuation
 - `getTopWords()`: Extract most frequent words
 - `saveResults()`: Export results to file

2. Benchmarking Infrastructure:

- **Dataset Generator:** Creates realistic text files with Zipf-like word distribution
- **Benchmark Runner:** Executes multiple runs, collects statistics
- **Performance Analyzer:** Generates graphs and summary reports

3. Workflow:

1. Generate test datasets (10MB - 100MB)
2. Compile sequential version with optimizations
3. Run 5 iterations per dataset
4. Collect timing statistics (mean, median, std dev)
5. Analyze bottlenecks and identify parallelization targets
6. Generate performance visualizations

6.3 Document Formatting and Compilation

6.3.1 LaTeX Formatting and Compilation

This document is written in **LaTeX** using Overleaf.

6.4 Project Implementation Summary

6.4.1 Parallel Implementation Approach

The parallel implementation successfully employs a **chunk-based parallelization strategy with thread-local storage**, as detailed in Section 3. The implementation uses OpenMP 5.x with the following key techniques:

Implemented Strategy: Thread-Local Maps with Reduction

```

1 #pragma omp parallel reduction(+ : totalWordCount)
2 {
3     WordMap localMap; // Thread-local hash map
4
5     #pragma omp for schedule(static)
6     for (int i = 0; i < rawWords.size(); ++i) {
7         std::string normalized = normalizeWord(rawWords[i]);
8         if (!normalized.empty()) {
9             localMap[normalized]++; // No sync needed
10            totalWordCount++; // Reduction aggregation
11        }
12    }
13
14    // Merge thread-local maps into global result
15    #pragma omp critical
16    {
17        for (const auto& entry : localMap) {
18            wordFreq[entry.first] += entry.second;
19        }
20    }
21 }

```

Listing 13: Implemented OpenMP Parallel Word Counter

Key Implementation Features:

- **Thread-Local Storage:** Each thread maintains private hash map
- **Zero Contention:** No synchronization during word counting phase
- **Reduction Clause:** Optimal counter aggregation (best performance)
- **Critical Merge:** Synchronized final combination of results
- **Static Scheduling:** Predictable load distribution for uniform workload

Performance Achievements:

- Best speedup: $1.81\times$ (10MB dataset, 8 threads)
- Reduction synchronization: 105% faster than critical sections
- Amdahl's Law validation: theoretical max $1.82\times$, achieved $1.81\times$
- 240 comprehensive benchmarks (4 datasets \times 3 sync methods \times 4 threads \times 5 runs)

For complete implementation details, see Section 3 (Parallel Implementation) and Section 4 (Performance Results).

7 References

1. OpenMP Architecture Review Board. *OpenMP Application Programming Interface Version 5.2*. November 2021. <https://www.openmp.org/specifications/>
2. Chapman, Barbara, Gabriele Jost, and Ruud van der Pas. *Using OpenMP: Portable Shared Memory Parallel Programming*. MIT Press, 2007.
3. Herlihy, Maurice, and Nir Shavit. *The Art of Multiprocessor Programming*. Morgan Kaufmann, 2012.
4. Pacheco, Peter. *An Introduction to Parallel Programming*. Morgan Kaufmann, 2011.
5. Amdahl, Gene M. *Validity of the single processor approach to achieving large scale computing capabilities*. AFIPS Conference Proceedings, 1967.
6. Williams, Anthony. *C++ Concurrency in Action, 2nd Edition*. Manning Publications, 2019.
7. Dr. Yasir Alguwaifli. *ARTI503 Course Materials - Parallel Computer Architecture*. IAU, 2025.
8. Amazon Web Services. *AWS Cloud Credits for Research*. <https://aws.amazon.com/research-credits/>

8 Appendix

8.1 Appendix A: Project Repository Structure

Complete directory tree:

```
parallel-grepper/
  src/
    sequential/          # Sequential implementation
      main.cpp
      word_counter_sequential.cpp
      word_counter_sequential.h
    parallel/            # Parallel implementation (OpenMP)
      main.cpp
      word_counter_parallel.cpp
      word_counter_parallel.h
  benchmarks/
    F-run_sequential_benchmarks.py # Sequential benchmarks
    F-run_parallel_benchmarks.py   # Parallel benchmarks
    F-analyze_sequential_results.py # Visualization & analysis
    F-analyze_parallel_results.py  # Visualization & analysis
    generate_dataset.py             # Test data generation
    results/                       # Benchmark outputs (CSV/JSON/PNG)
  scripts/
    build.ps1                    # Windows build automation
    build.sh                     # Linux/macOS build script
    run_sync_tests.ps1          # Race condition tests
  docs/
    proposal.tex                 # This document
    BUILD_GUIDE.md               # Compilation instructions
    RACE_CONDITION_FIXES.md      # Synchronization method details
  data/                          # Test datasets
    # Generated text files (10MB - 100MB+) from generate_dataset.py
  results/
    sequential/                  # Sequential outputs
    parallel/                     # Parallel outputs
```

8.2 Appendix B: Build and Run Instructions

Prerequisites:

- GCC 15.2.0+ with OpenMP support (MinGW-w64 POSIX UCRT on Windows)
- C++17 standard library
- Python 3.8+ with pandas, matplotlib, seaborn (for benchmarking)

Compilation Commands:

```
1 # Create build directory
2 mkdir -p build results/sequential
3
4 # Compile sequential version (Linux/macOS)
5 g++ -std=c++17 -O3 -march=native \
6     -o build/sequential_counter \
7     src/sequential/word_counter_sequential.cpp \
8     src/sequential/main.cpp
9
10 # Windows with MinGW
11 g++ -std=c++17 -O3 ^
12     -o build/sequential_counter.exe ^
13     src/sequential/word_counter_sequential.cpp ^
14     src/sequential/main.cpp
```

Listing 14: Sequential Version Build

```
1 # Create build directory
2 mkdir -p build results/parallel
3
4 # Compile parallel version (Linux/macOS)
5 g++ -std=c++17 -O3 -fopenmp -march=native \
6     -o build/parallel_counter \
7     src/parallel/word_counter_parallel.cpp \
8     src/parallel/main.cpp
9
10 # Windows with MinGW
11 g++ -std=c++17 -O3 -fopenmp ^
12     -o build/parallel_counter.exe ^
13     src/parallel/word_counter_parallel.cpp ^
14     src/parallel/main.cpp
```

Listing 15: Parallel Version Build

Usage Examples:

```

1 # Default: 8 threads with reduction synchronization
2 ./build/parallel_counter data/test_100mb.txt
3
4 # Specify thread count
5 ./build/parallel_counter data/test_50mb.txt 4
6
7 # Specify sync method: reduction, atomic, or critical
8 ./build/parallel_counter data/test_10mb.txt 4 reduction
9 ./build/parallel_counter data/test_10mb.txt 4 atomic
10 ./build/parallel_counter data/test_10mb.txt 4 critical
11
12 # Windows PowerShell
13 .\build\parallel_counter.exe data\test_100mb.txt 8 reduction

```

Listing 16: Running Parallel Counter

8.3 Appendix C: Performance Metrics Definitions**Key Performance Indicators:****1. Speedup:**

$$S(p) = \frac{T_{\text{sequential}}}{T_{\text{parallel}}(p)}$$

where p is the number of threads

2. Efficiency:

$$E(p) = \frac{S(p)}{p} \times 100\%$$

Ideal efficiency is 100%, indicating perfect scaling

3. Throughput:

$$\text{Throughput} = \frac{\text{Data Size (MB)}}{\text{Execution Time (seconds)}}$$

4. Strong Scalability:

Fixed problem size, varying thread count

Measured by how speedup changes as threads increase

5. Amdahl's Law:

$$S_{\text{max}}(p) = \frac{1}{(1 - P) + \frac{P}{p}}$$

where P is the parallelizable fraction

8.4 Appendix D: Synchronization Method Code Examples

Reduction (Recommended):

```

1 #pragma omp parallel reduction(+: totalWordCount)
2 {
3     WordMap localMap;
4     #pragma omp for schedule(static)
5     for (int i = 0; i < rawWords.size(); ++i) {
6         std::string normalized = normalizeWord(rawWords[i]);
7         if (!normalized.empty()) {
8             localMap[normalized]++;
9             totalWordCount++; // Automatic reduction
10        }
11    }
12    // Merge localMap into wordFreq under critical section
13 }

```

Listing 17: Reduction Synchronization

Atomic Operations:

```

1 #pragma omp parallel
2 {
3     WordMap localMap;
4     #pragma omp for schedule(static)
5     for (int i = 0; i < rawWords.size(); ++i) {
6         std::string normalized = normalizeWord(rawWords[i]);
7         if (!normalized.empty()) {
8             localMap[normalized]++;
9             #pragma omp atomic
10            totalWordCount++; // Atomic increment
11        }
12    }
13    // Merge phase
14 }

```

Listing 18: Atomic Synchronization

Critical Sections:

```

1 #pragma omp parallel
2 {
3     WordMap localMap;
4     #pragma omp for schedule(static)
5     for (int i = 0; i < rawWords.size(); ++i) {
6         std::string normalized = normalizeWord(rawWords[i]);
7         if (!normalized.empty()) {
8             localMap[normalized]++;
9             #pragma omp critical
10            totalWordCount++; // Mutex-protected increment
11        }
12    }
13    // Merge phase
14 }

```

Listing 19: Critical Section Synchronization

8.5 Appendix E: Acknowledgments

This project was completed as part of the **Parallel Computer Architecture (ARTI503)** course at **Imam Abdulrahman Bin Faisal University (IAU)**, College of Computer Science and Information Technology (CCSIT).

Course Details:

- Fourth Year, Semester 1, Academic Year 2025-2026
- Instructor: Dr. Yasir Alguwaifli
- Project: Parallel Word Frequency Counter
- Framework: OpenMP 5.x for shared-memory parallelization

Tools and Technologies:

- C++17 with GCC 15.2.0 (MinGW-w64 POSIX UCRT)
- OpenMP 5.x API for parallel programming
- Python 3.x with pandas, matplotlib, seaborn for analysis
- Windows 11 development environment
- LaTeX/Overleaf for documentation