

2024 Fall CMPT 431 Project Report Fast Fourier Transform

Henry Chau, Bogdan Badea

Introduction (what is the problem you're solving)

The Fast Fourier Transform (FFT) is a computationally efficient algorithm for calculating the Discrete Fourier Transform (DFT), which transforms signals from their original domain (usually time or space) to the frequency domain. It is one of the most important algorithms in science and engineering, with numerous real-world applications, including signal analysis, audio filtering, and image compression.

In this project, we will implement three versions of the FFT:

- A serial version (single-threaded, single-process)
- A parallel version (single-threaded, single-process)
- A distributed version (single/multi-threaded, multi-process)

We will test and compare these implementations to find out which one is the most efficient. We will also look into what makes the fastest FFT algorithms work better and how they improve performance.

Background (what is the state of the art in serial and parallel implementations of the problem)

The basic DFT is shown in the equation below. When calculated directly, it has a time complexity of $O(n^2)$, which can be slow for large data sets.

$$Y[i] = \sum_{k=0}^{n-1} X[k] \omega^{ki}, \quad 0 \leq i < n, \quad \omega = e^{2\pi\sqrt{-1}/n}$$

The FFT is a divide-and-conquer algorithm for computing the DFT, introduced by James Cooley and John Tukey in 1965[1]. Originally developed to detect nuclear tests by the Soviet Union, it later proved to be very useful in many other applications. The FFT works by recursively breaking down the original data into two smaller parts, which reduces the time complexity to $O(n \log n)$.

The equation below shows how the DFT is split into smaller problems(odd and even parts). There are different types of FFT algorithms, such as the Prime-factor FFT, but this project focuses on the Radix 2 Cooley-Tukey FFT.

$$Y[i] = \sum_{k=0}^{(n/2)-1} X[2k]\omega^{2ki} + \sum_{k=0}^{(n/2)-1} X[2k+1]\omega^{(2k+1)i}$$

Fastest Fourier Transform in the West(FFTW) is a popular and efficient implementation of the Fast Fourier Transform, developed by Matteo Frigo and Steven Johnson. It supports serial, parallel, and distributed versions of FFT. FFTW is highly optimized, customizable, and can adapt to different hardware configurations. It uses a plan-based approach: before computation begins, a plan is created at runtime to select the best combination of codelets (optimized C code blocks generated by a specialized compiler)[2]. The executor then runs these codelets to perform the transformation.

Another option is to use the parallel processing power of GPUs to speed up the algorithm. Nvidia's CUDA Fast Fourier Transform library (cuFFT) is a popular tool for performing FFT efficiently on GPUs.

Implementation Details (what you have implemented in the serial, parallel and distributed versions)

For simplicity, we restrict the input size n to powers of two. Inputs that are not powers of two can still be used by adding zero padding.

FFT and inverse FFT (IFFT) are nearly identical operations, with two key differences: the exponent sign of ω is positive in IFFT, and the IFFT results must be normalized at the end. As a result, both can be combined into a single function. A boolean argument is used to determine the direction.

Serial

This is a recursive version of the serial implementation. While it is simple and easy to read, its recursive nature makes it hard to optimize and adapt for parallel version.

```

1  std::vector<std::complex<double>> recursiveFFT(
2      const std::vector<double> &X, std::complex<double> omega) {
3      std::vector<std::complex<double>> Y(X.size());
4      const int n = X.size();
5      if (n == 1) {
6          Y[0] = X[0];
7          return Y;
8      }
9      // Populate new vectors containing only odd or even indexed elements
10     std::vector<double> xEven;
11     for (int i = 0; i < n; i += 2) {
12         xEven.push_back(X[i]);

```

```

13     }
14     std::vector<double> xOdd;
15     for (int i = 1; i < n; i += 2) {
16         xOdd.push_back(X[i]);
17     }
18     std::vector<std::complex<double>> left = recursiveFFT(xEven, std::pow(omega, 2));
19     std::vector<std::complex<double>> right = recursiveFFT(xOdd, std::pow(omega, 2));
20     for (int i = 0; i < n; ++i) {
21         Y[i] = left[i % (n / 2)] + std::pow(omega, i) * right[i % (n / 2)];
22     }
23     return Y;
24 }

```

This is the FFT algorithm proposed in the *Introduction to Parallel Computing(IPC)* textbook[3]. During computation, the input array remains unordered, requiring bit manipulation to determine the correct indices of j , k and the twiddle factor. A cache is used to avoid repeated calculation of omega. The main advantage of this algorithm is that the results computed at each stage are stored contiguously, making it well-suited for data exchange in multi-process setups. The result needs to be reordered at the end of the algorithm.

```

1  std::vector<std::complex<double>> iterativeIcpFft(
2      const std::vector<std::complex<double>> &X, bool isInverse) {
3      const int n = X.size();
4      const int r = std::log2(n);
5      const std::complex<double> img(0.0, 1.0);
6      double exponentSign = isInverse ? 1.0 : -1.0;
7      std::complex<double> omegaBase = std::exp(exponentSign * 2.0 * img
8          * std::numbers::pi / (double)n);
9      // Result array
10     std::vector<std::complex<double>> R(X);
11     // Auxillary array to hold previous value of R
12     std::vector<std::complex<double>> S(R);
13     // Store previously calculated omega
14     std::map<int, std::complex<double>> omegaCache;
15     // Outer loop  $O(\log n)$ ,  $m$  represent stage
16     for (int m = 0; m < r; ++m) {
17         // Inner loop  $O(n)$ 
18         for (int i = 0; i < n; ++i) {
19             // Let  $(b_0 b_1 \dots b_{r-1})$  be the binary representation of  $i$ 
20             // Clear the bit,  $j := (b_0 \dots b_{m-1} 0 b_{m+1} \dots b_{r-1})$ 
21             unsigned int j = i & ~(1 << (r - 1 - m));
22             // Set the bit,  $k := (b_0 \dots b_{m-1} 1 b_{m+1} \dots b_{r-1})$ 
23             unsigned int k = i | (1 << (r - 1 - m));
24             //  $(b_m b_{m-1} \dots b_0 \dots 0 \dots 0)$ 
25             unsigned int omegaExp = getFirstNBits(reverseBits(i, r), m + 1) << (r - 1 - m);
26             std::complex<double> omega;
27             if (omegaCache.find(omegaExp) != omegaCache.end()) {
28                 omega = omegaCache[omegaExp];
29             } else {
30                 omega = std::pow(omegaBase, omegaExp);
31                 omegaCache[omegaExp] = omega;
32             }
33             R[i] = S[j] + S[k] * omega;

```

```

34     }
35     // Update S with new values
36     std::copy(R.begin(), R.end(), S.begin());
37 }
38 // In-place bit-reversal reordering
39 for (int i = 0; i < n; ++i) {
40     int reversedIndex = reverseBits(i, r);
41     if (i < reversedIndex) {
42         std::swap(R[i], R[reversedIndex]);
43     }
44 }
45 // Normalize result if performing IFFT
46 if (isInverse) {
47     for (int i = 0; i < n; ++i) {
48         R[i] /= n;
49     }
50 }
51 return R;
52 }

```

Parallel

From the iterativeIcp algorithm we can note that each of the iterations of the outerloop depends on the results of the previous iteration, but that there is no such dependency for the inner loop, and that the inner loops only modify their respective indices of the result array, meaning we don't any mutexes for the writes, only barrier to make sure they all stay on the same iteration, and an atomic value to allow them to dynamically go through the array since some indices have more work then others as a result of the twiddlers or index values.

```

1  std::atomic<int> iterativeIcpFft_covered_i(0);
2  void multithreaded_iterativeIcpFft_loop_function(const int r, const int n,
3          std::vector<std::complex<double>> &R,
4          std::vector<std::complex<double>> &S,
5          const std::complex<double> omega, std::barrier<> &b, int thread_id,
6          int block_size) {
7      int start_i;
8      int m;
9      int i;
10     unsigned int omegaExp;
11     // Outer loop O(log n)
12     for (m = 0; m < r; ++m) {
13         if (thread_id == 0) {
14             std::copy(R.begin(), R.end(), S.begin());
15             iterativeIcpFft_covered_i = 0;
16         }
17         b.arrive_and_wait();
18         // Inner loop O(n)
19         start_i = iterativeIcpFft_covered_i.fetch_add(block_size);
20         while (start_i < n) {
21             for (i = start_i; i < std::min(start_i + block_size, n); ++i) {
22                 // Let (b_0 b_1 ... b_{r-1}) be the binary representation of i
23                 std::bitset<32> j(i);

```

```

24         j.reset(r - 1 - m); // (b_0 ... b_m1 0 b_m+1 ... b_r1)
25         std::bitset<32> k(i);
26         k.set(r - 1 - m); // (b_0 ... b_m1 1 b_m+1 ... b_r1)
27
28         omegaExp = getFirstNBits(reverseBits(i, r), m + 1) << (r - 1 - m);
29         // (b_m b_m-1 ... b-0 ... 0 ... 0)
30         R[i] = S[j.to_ulong()] + S[k.to_ulong()] * std::pow(omega, omegaExp);
31     }
32     start_i = iterativeIcpFft_covered_i.fetch_add(block_size);
33 }
34 b.arrive_and_wait();
35 }
36 }
37
38 std::vector<std::complex<double>> multithreaded_iterativeIcpFft(
39     const std::vector<std::complex<double>> &X,
40     bool isInverse, int nThreads, int block_size) {
41     const int n = X.size();
42     const int r = std::log2(n);
43     std::barrier b(nThreads);
44     const std::complex<double> img(0.0, 1.0);
45     double exponentSign = isInverse ? 1.0 : -1.0;
46     std::complex<double> omega = std::exp(exponentSign * 2.0 * img * std::numbers::pi / (double)n);
47     std::vector<std::complex<double>> R(X); // Result array
48     std::vector<std::complex<double>> S(n); // Auxillary array to hold previous value of R
49     std::vector<std::thread> threads;
50     for (int i = 0; i < nThreads; ++i) {
51         threads.push_back(std::thread(multithreaded_iterativeIcpFft_loop_function, r, n, std::ref(R),
52             std::ref(S), omega, std::ref(b), i, block_size));
53     }
54     for (auto &t : threads) {
55         t.join();
56     }
57     // In-place bit-reversal reordering
58     for (int i = 0; i < n; ++i) {
59         int reversedIndex = reverseBits(i, r);
60         if (i < reversedIndex) {
61             std::swap(R[i], R[reversedIndex]);
62         }
63     }
64     // Normalize result if performing IFFT
65     if (isInverse) {
66         for (int i = 0; i < n; ++i) {
67             R[i] /= n;
68         }
69     }
70     return R;
71 }

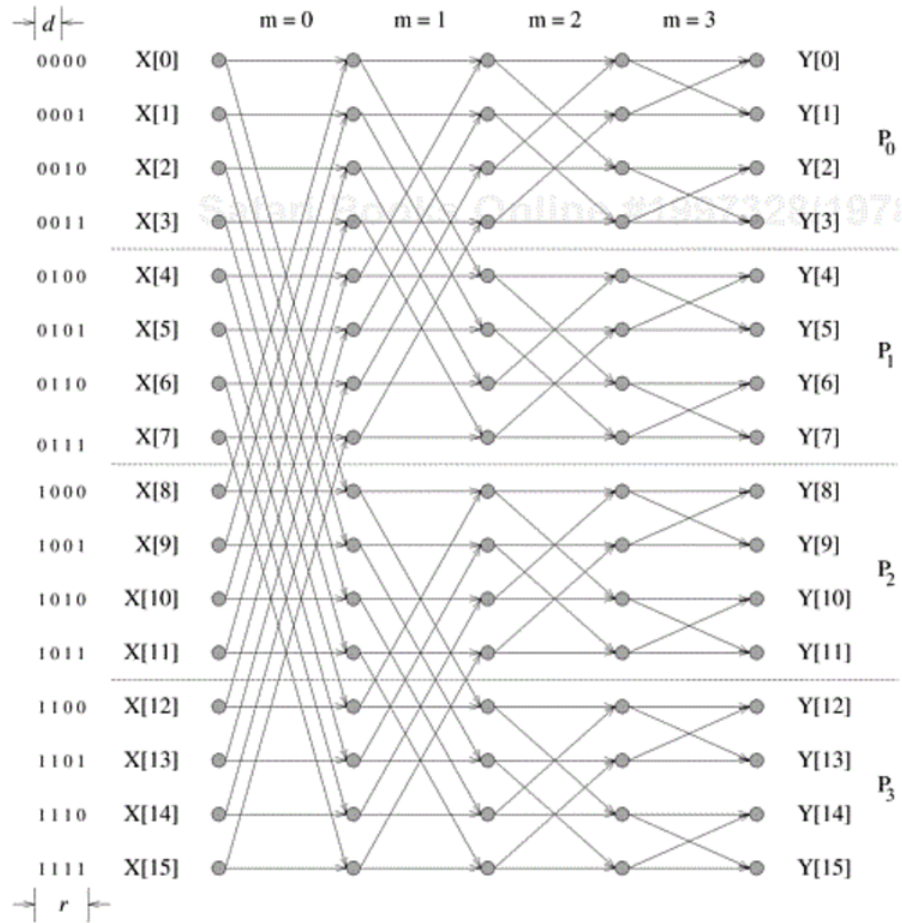
```

Distributed

The figure below, taken from *IPC Ch.13*[3], illustrates how process communication works between each stage in the binary exchange algorithm. In this approach, we observe that $p/2$ pairs of processes

communicate with each other during each communication stage, with each process receiving a new partner at every stage. The communication is performed using a blocking *MPI_Sendrecv* to ensure that the send and receive operations occur almost simultaneously, eliminating the risk of deadlock. Notably, communication only occurs during stages where $m < \log p$

For the last $\log n - \log p$ iterations of the outer loop, the computation is entirely local. This means that these iterations can be treated as a "black-box" operation, making it feasible to utilize other libraries or implement multithreading to further accelerate the local computation.



```

1  std::vector<std::complex<double>> icpFftDistributed(
2      const std::vector<std::complex<double>> &X, bool isInverse) {
3      int processes, currProcId;
4      MPI_Comm_size(MPI_COMM_WORLD, &processes); // Total number of processes
5      MPI_Comm_rank(MPI_COMM_WORLD, &currProcId); // Current process ID
6      const int n = X.size();
7      const int r = std::log2(n);
8      const int logP = std::log2(processes);
9      const int elePerProc = n / processes;
10     // Start index of elements owned by this process
11     const int startIdx = currProcId * elePerProc;
12     const int endIdx = startIdx + elePerProc; // Exclusive end index
13     const std::complex<double> img(0.0, 1.0);
14     double exponentSign = isInverse ? 1.0 : -1.0;

```

```

15 std::complex<double> omegaBase = std::exp(exponentSign * 2.0 * img
16      * std::numbers::pi / (double)n);
17 std::vector<std::complex<double>> R(n); // Result array
18 std::copy(X.begin() + startIdx, X.begin() + endIdx, R.begin() + startIdx);
19 std::vector<std::complex<double>> S(n); // Auxillary array to hold previous value of R
20 std::copy(R.begin() + startIdx, R.begin() + endIdx, S.begin() + startIdx);
21 std::unordered_map<int, std::complex<double>> omegaCache; // Store previously calculated omega
22 // Outer loop O(log n), m represent stage
23 for (int m = 0; m < r; ++m) {
24     // Check if data exchange is needed at this stage
25     if (m < (logP)) {
26         std::bitset<32> partnerProcIdBits(currProcId);
27         partnerProcIdBits.flip(logP - m - 1);
28         int partnerProcId = partnerProcIdBits.to_ulong();
29         // Exchange data with partner
30         MPI_Sendrecv(
31             R.data() + startIdx, elePerProc, MPI_DOUBLE_COMPLEX, partnerProcId, m,
32             S.data() + partnerProcId * elePerProc, elePerProc, MPI_DOUBLE_COMPLEX,
33             partnerProcId, m, MPI_COMM_WORLD, MPI_STATUS_IGNORE
34         );
35     }
36     // Inner loop O(n)
37     for (int i = startIdx; i < endIdx; ++i) {
38         // Let (b_0 b_1 ... b_{r-1}) be the binary representation of i
39         // Clear the bit, j := (b_0 ... b_{m-1} 0 b_{m+1} ... b_{r-1})
40         unsigned int j = i & ~(1 << (r - 1 - m));
41         // Set the bit, k := (b_0 ... b_{m-1} 1 b_{m+1} ... b_{r-1})
42         unsigned int k = i | (1 << (r - 1 - m));
43         // (b_m b_{m-1} ... b_0 ... 0 ... 0)
44         unsigned int omegaExp = getFirstNBits(reverseBits(i, r), m + 1) << (r - 1 - m);
45         std::complex<double> omega;
46         if (omegaCache.find(omegaExp) != omegaCache.end()) {
47             omega = omegaCache[omegaExp];
48         } else {
49             omega = std::pow(omegaBase, omegaExp);
50             omegaCache[omegaExp] = omega;
51         }
52         R[i] = S[j] + S[k] * omega;
53     }
54     // Update S with new values
55     std::copy(R.begin() + startIdx, R.begin() + endIdx, S.begin() + startIdx);
56 }
57 // All processes gather result from other processes
58 MPI_Allgather(R.data() + startIdx, elePerProc, MPI_DOUBLE_COMPLEX, R.data(),
59     elePerProc, MPI_DOUBLE_COMPLEX, MPI_COMM_WORLD);
60 // In-place bit-reversal reordering
61 for (int i = 0; i < n; ++i) {
62     int reversedIndex = reverseBits(i, r);
63     if (i < reversedIndex) {
64         std::swap(R[i], R[reversedIndex]);
65     }
66 }
67 // Normalize result if performing IFFT
68 if (isInverse) {

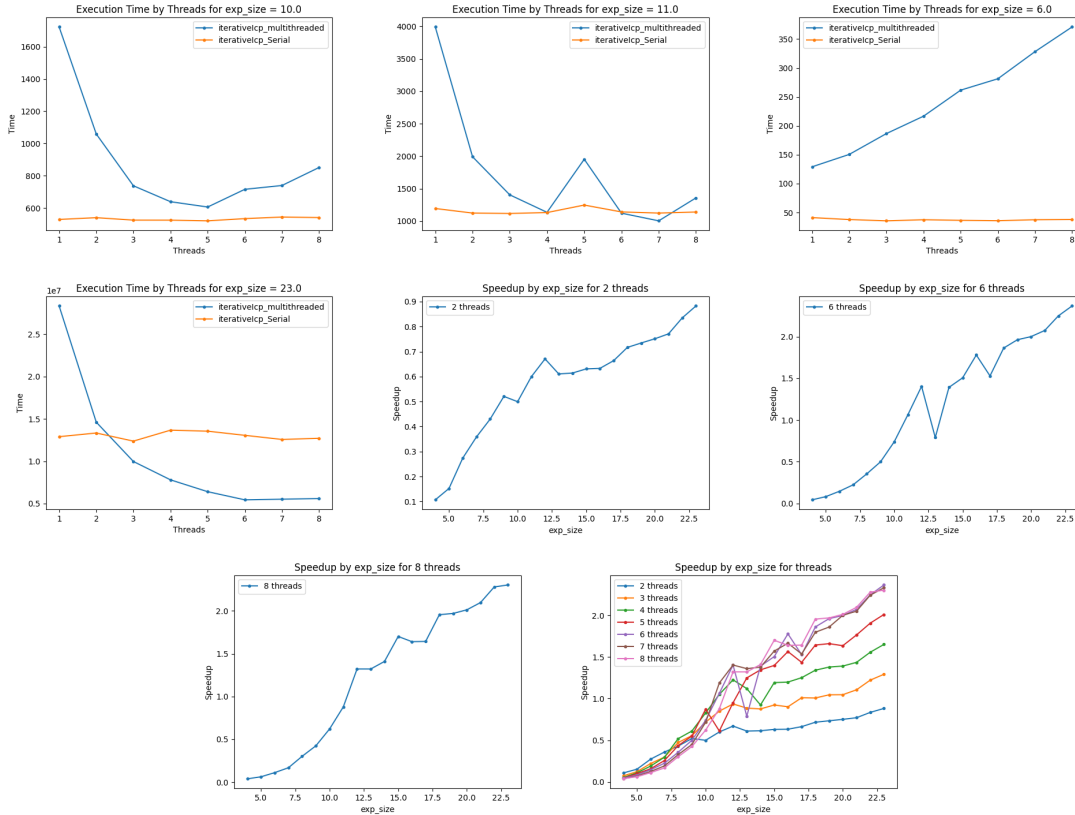
```

```

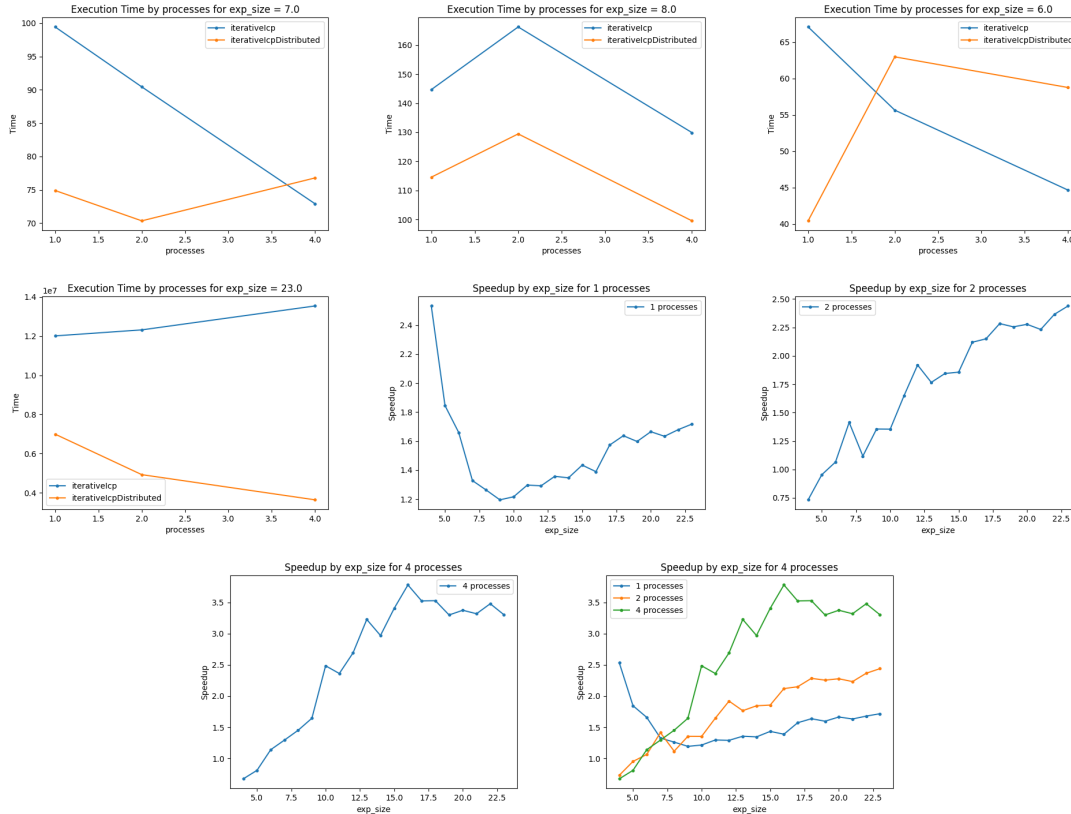
69     for (int i = 0; i < n; ++i) {
70         R[i] /= n;
71     }
72 }
73 return R;
74 }

```

Evaluation (results from your experiments)



In our exploration of the parallel implementation, we discovered that leveraging multiple threads only became advantageous when the problem size reached approximately $\sim 2^{11}$. Before this point, the added overhead from thread management outweighed any computational benefits, leading to longer execution times rather than speedups. This suggests that for smaller problem sizes, the cost of parallelization can be prohibitive. However, as the exponent x in 2^x grew larger, the performance gains from parallel processing became increasingly evident. The speedup improved consistently across all thread counts for sufficiently large problem sizes. Notably, the performance was most effective with higher thread counts, reaching optimal results around six threads, beyond which the benefits began to plateau.



In the case of the distributed version, we found that it incurred significantly less overhead compared to the parallel implementation. This reduced penalty allowed it to become effective at smaller problem sizes, with noticeable benefits starting at around 2^7 rather than 2^{11} . Additionally, the speedup achieved by the distributed version was impressive, closely approaching the ideal of perfect parallelization when a small number of processes were used. However, as the number of processes increased, the speedup began to deviate more noticeably from this ideal, with its proportional gains diminishing most visibly for 4 processes and high exp_size . This highlights the efficiency of the distributed approach for smaller numbers of processes while also illustrating the beginning of diminishing returns that occur as more processes are introduced.

Conclusions (summary of results/speedups, what you learned from the project)

Parallelizing the Fast Fourier Transform (FFT) proves to be significantly more efficient when utilizing OpenMPI for distributed processing compared to multithreading. This is evident from our results, which show that even with just four processes, the distributed implementation surpasses the upper limit of the speedup achieved by six threads in the multithreaded version. This suggests that distributed processing is better suited for scaling the FFT, as it manages to overcome the limitations imposed by thread-based approaches, particularly for larger problem sizes. While parallelizing the FFT can offer substantial performance improvements, it is worth noting that the most effective strategy may involve leveraging specialized hardware and libraries designed

for these tasks. For instance, NVIDIA's CUDA Fast Fourier Transform library (cuFFT) offers highly optimized FFT implementations that are specifically designed to run on GPUs. Offloading computational workloads to a GPU using cuFFT can provide far greater performance gains than either threading or distributed processing on CPUs[4]. This highlights the importance of selecting the right computational framework and hardware platform when optimizing FFT computations, especially for applications requiring high-performance solutions.

References

- [1] James W. Cooley, John W. Tukey(1965), An Algorithm for the Machine Calculation of Complex Fourier Series
- [2] Matteo Frigo, Steven G. Johnson(1998), FFTW: An Adaptive Software Architecture for the FFT
- [3] Ananth Grama (2003), Introduction to Parallel Computing, Addison-Wesley
- [4] C. Zhang, Y. Xiao, H. Chen and H. Liu, "Parallel Optimization and Hardware Customization for Fast Fourier Transform," 2023 IEEE 14th International Conference on Software Engineering and Service Science (ICSESS), Beijing, China, 2023, pp. 105-110