

Algorithms Implemented

Three different approaches were used to find the sum of the arrays. The first approach is the binary tree PRAM summation. For this approach, the array is divided amongst the number of threads, where each thread computes a set number of two element sums and then places them into an array of size $N/2$. The function then recurses, passing down this array of $N/2$ to be summed. Although we know the theoretical parallel time complexity to be $O(\log_2 N)$, looking at the serial complexity, the following can be observed:

$$O(\text{binary_sum}) = O\left(N + \frac{N}{2} + \frac{N}{4} + \cdots + \frac{N}{2^i} + \cdots + 2\right) = O(2N - 1) = O(N)$$

The second approach was to divide the array across the number of processors and sum all elements assigned to that processor, place that partial sum into a scratch array with the size equal to the number of threads, and then sum all of the scratch array elements afterwards. This divide-and-conquer function was named `parallel_sum`. When the number of threads is equal to 1, the order

$$\begin{aligned} O(\text{parallel_sum}) &= O(N); p = 1 \\ O(\text{parallel_sum}) &= O(N/p) \end{aligned}$$

The third approach is, in theory, identical to the parallel sum approach. By implementing the parallel for reduction method of OpenMP, the array should be split into equal segments per thread, each section added into a temporary sum, and then the actual sum computed by summing all of these partial sums. The time complexity should be the same as that of the previous approach. In practice, however, we will observe that the extra overhead in manually managing the divide-and-conquer method causes the time to be slightly higher than that of the parallel reduction method.

All three approaches were run using 1, 2, 4, 8, 16, 20, and 24 threads. One preliminary result that was interesting is that, even for the divide-and-conquer methods, as the number of threads diverged from powers of two performance was reduced. Using 16 threads had the best performance for all algorithms.

Algorithm Performance Comparison

For all of the graphs in this section, time is on the Y axis, the number of array elements is on the X axis and the number of threads is indicated by the legend. This first graph plots the performance of the binary tree summation. For all numbers of threads, the amount of time spiked at $N = 8,200,000$. As stated before, each power of two increases the performance (which is to be expected), however using 20 threads has a lower performance than even 8 threads and using 24 threads is only slightly better than 8 threads.

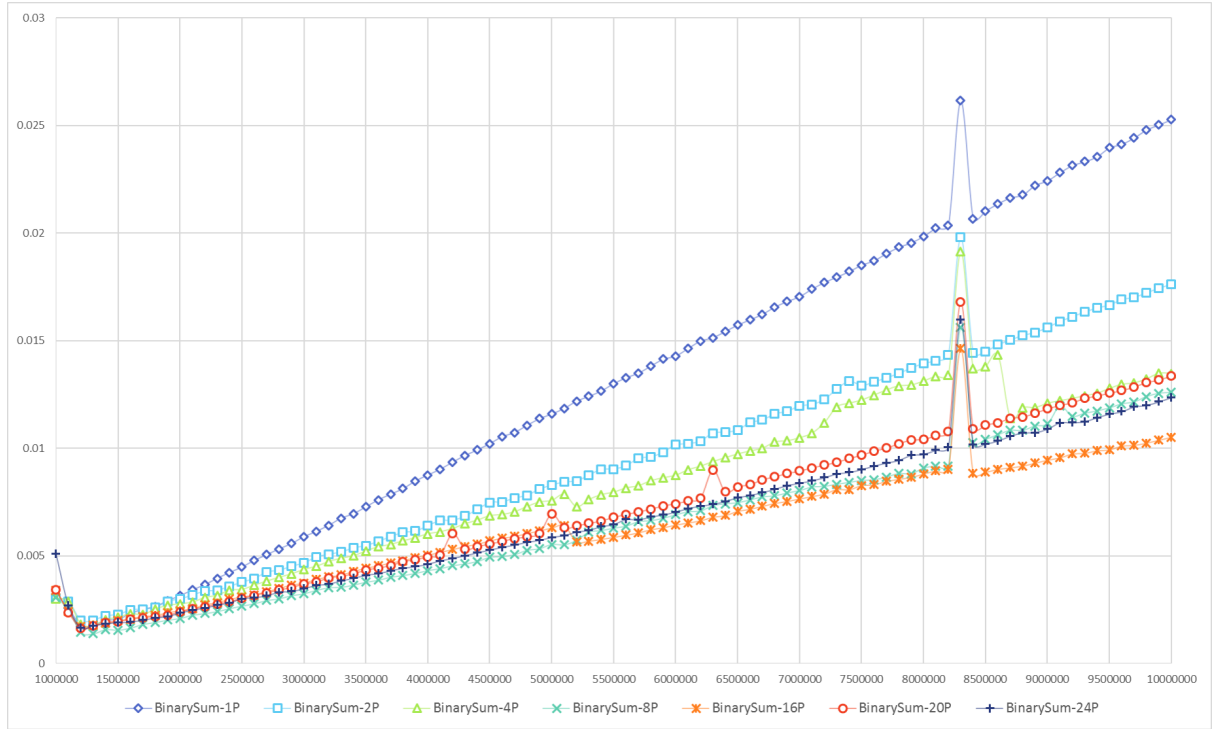


Figure 1: Binary Tree Sum Performance

Second up is the parallel sum method. Even in the worst case (1 thread) it is still more than twice as fast as the binary sum method. This follows our realistic expectations of the asymptotic complexity. The time complexity of the serial binary sum is $O(2N - 1)$ while the parallel sum method only requires $O(N)$ time. As the number of threads increases, this difference becomes more and more drastic. The same performance spike around 8.2 million is present here as well. It is also worth noting that there is odd behavior when using 4 threads.

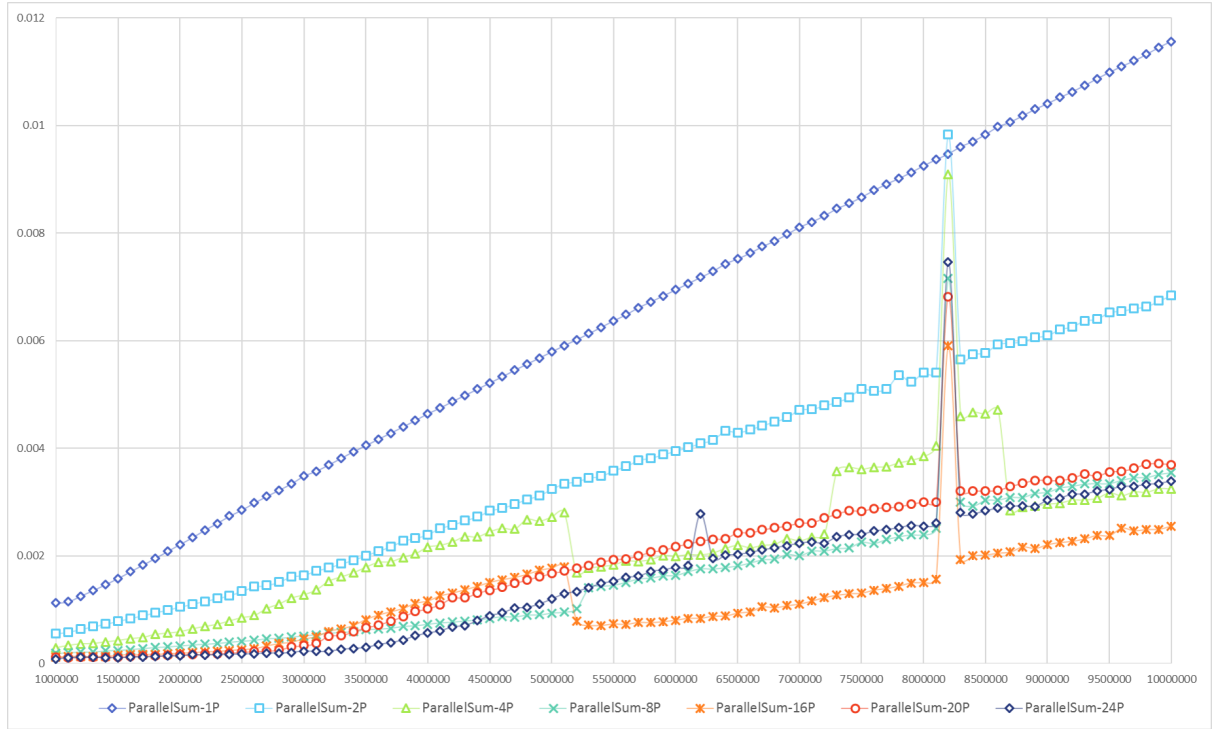


Figure 2: divide-and-conquer Sum Performance

When observing the fast sum method, we see very similar values of time for threads and, as with the previous method, there is odd stepping in performance when using 4 threads. It is worth noting, however, that while in the previous two methods using 20 threads was worse than 8 threads, when using the reduction method, the performance is almost equivalent. The major feature that stands out is that using this method does not cause the peak in performance around 8.2 million elements.

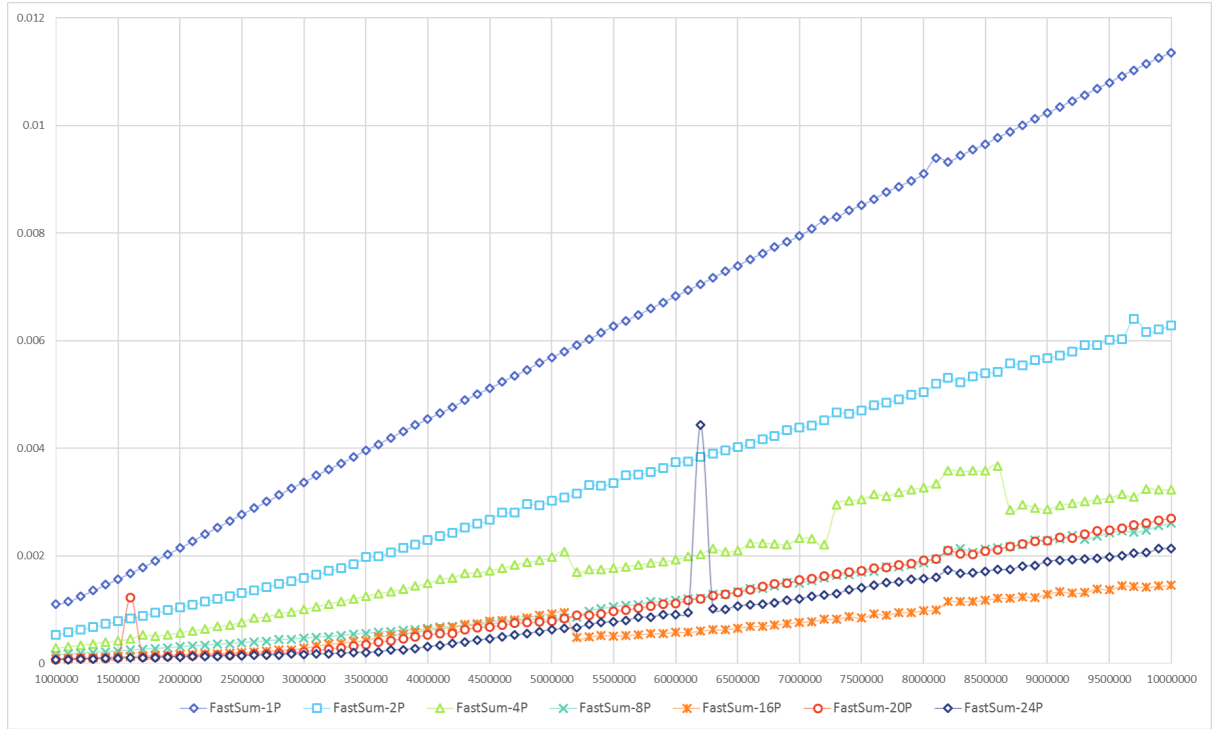


Figure 3: Reduction Sum Performance

As a side-by-side comparison of the two methods, the fast sum and parallel sum performances are plotted below. The fast sum values are marked by the dashed line. As expected, by manually managing the shared variables and scratch space used by the parallel method, the performance is slightly worse on all accounts.

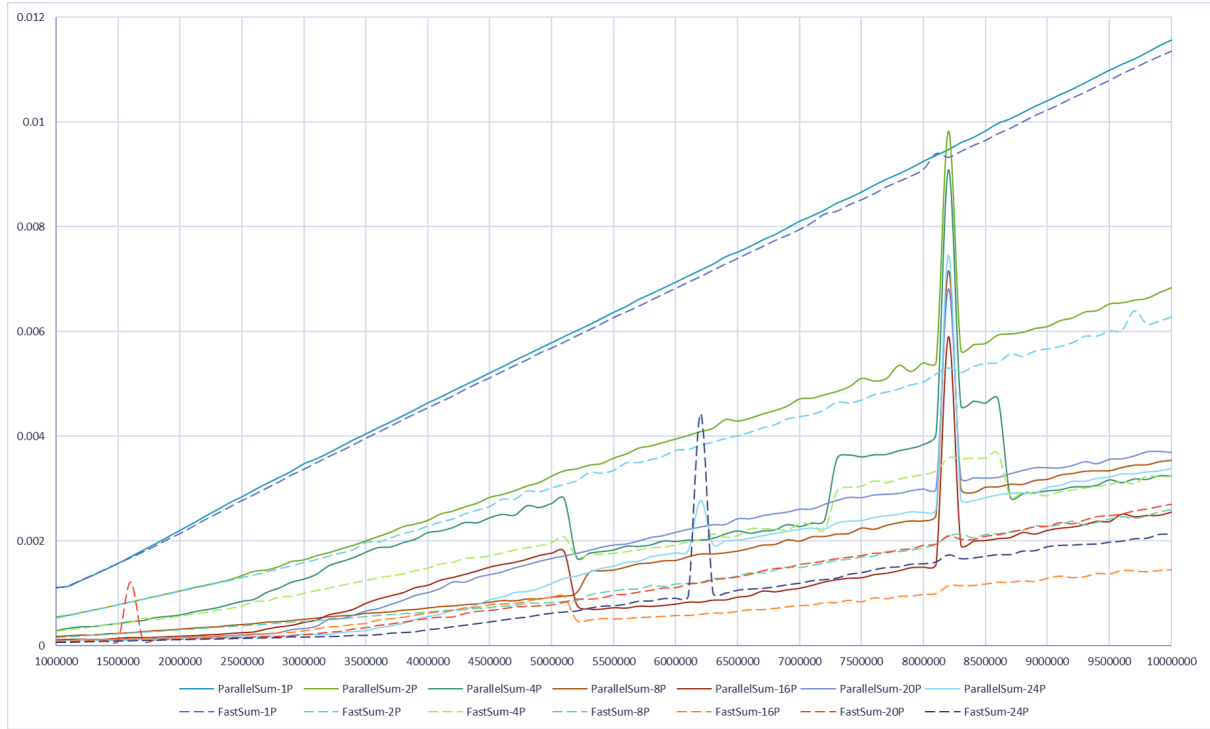


Figure 4: Parallel vs Fast Sum Performance

This graph provides a simple comparison between the serial performance of all three methods. Again, as assumed, the binary sum takes nearly twice the time that the parallel and fast methods do at large N . The management of the scratch space used also takes a toll on $N < 1,300,000$, as can be observed by the significantly higher runtime for these smaller N values in the binary tree sum.

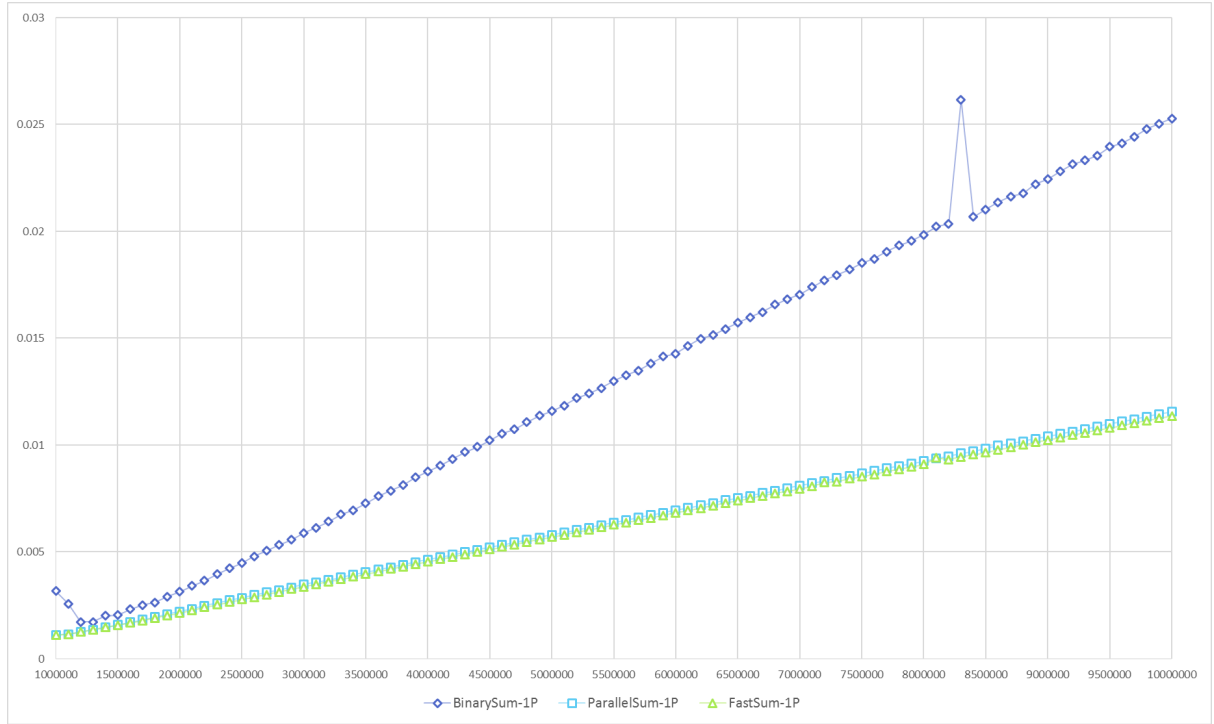


Figure 5: Single Thread Performance

The final chart shows the performance of each method when using 16 threads, as this number of threads produced the best performance. We can see, again that the management required for the binary sum creates a very large disparity between the runtimes when compared to the other two methods. At the largest value of N , the runtime of the binary sum is $\approx 3.64x$ that of the manually managed divide-and-conquer method and $\approx 5.79x$ of the fast sum runtime.

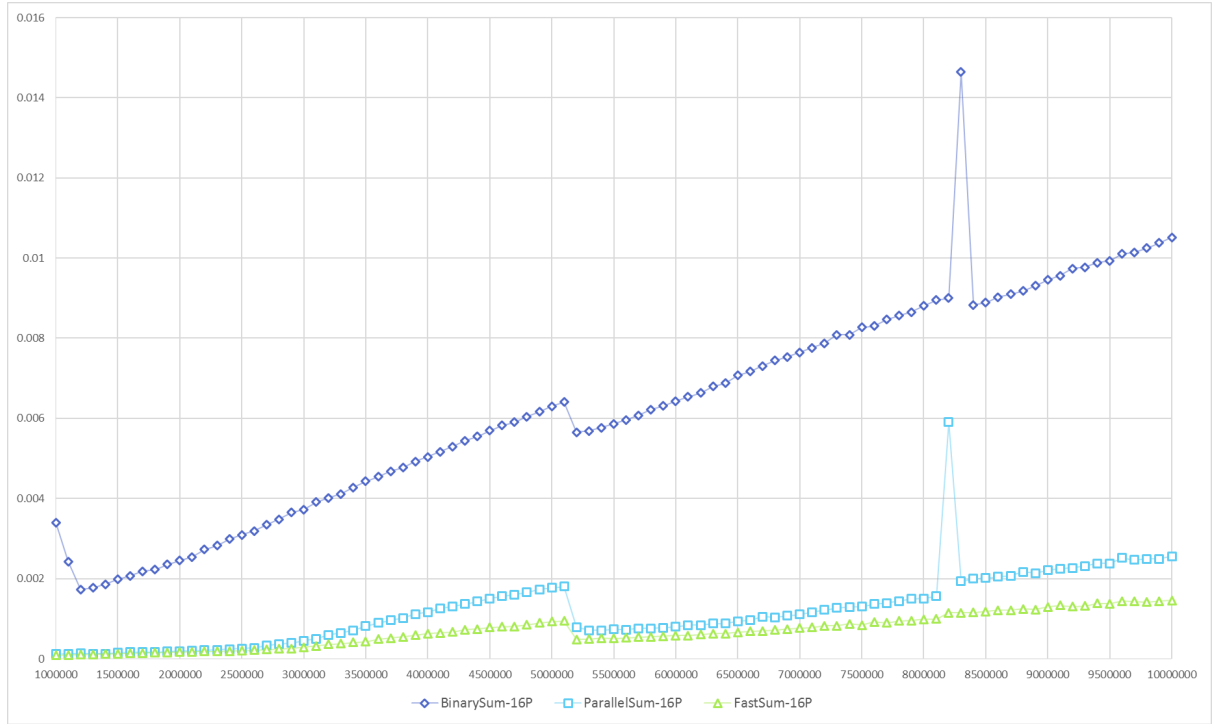


Figure 6: 16 Threads Performance

Conclusion

Although we know the theoretical limit of the binary sum to be $O(\log_2 N)$, in practice the actual parallel performance is closer to that of the serial performance. Even at small values of N , which we could realistically have more processors than data points, the standard parallel sum may actually run faster due to the amount of overhead it takes to create the scratch space for the binary tree sum method.

Code Appendix

```
1  /*****
2  Compile command: g++ -std=c++0x -fopenmp -o assignment3 assignment3.cpp -O3
3  *****/
4  #include <omp.h>
5  #include <iostream>
6  #include <cmath>
7  #include <random>
8  #include <sys/time.h>
9  #include <fstream>
10
11 void get_walltime(long double *wcTime)
12 {
13     struct timeval tp;
14     gettimeofday(&tp, NULL);
15     *wcTime = (long double) (tp.tv_sec + tp.tv_usec / 1000000.0);
16 }
17
18 /** binary_sum
19 *
20 * Calculates the sum of an array in a recursive binary-tree PRAM fashion.
21 *
22 * @param array - the array to be computed
23 * @param array_size - necessary for proper assignment
24 * @param time - reference to a variable to store the walltime
25 *
26 * @return - sum of array as a double
27 */
28 double binary_sum(double *array, int array_size, long double *time)
29 {
30     long double start_time = 0.0, end_time = 0.0, scratch_time = 0.0;
31     double answer = 0.0, *results;
32     int scratch_size, num_elements;
33
34     // Handle edge cases of 1 or 2 element arrays
35     if (array_size == 1)
36     {
37         return (array[0]);
38     }
39     if (array_size == 2)
40     {
41         return (array[0] + array[1]);
42     }
43
44     /** Create a scratch array to put temporary results in. If the input array is
45     not divisible by 2, add 1 so that
46     * scratch array can be a round integer. */
47     if (array_size % 2 == 1)
48     {
49         scratch_size = (array_size + 1) / 2;
50         results = new double[scratch_size]();
51     } else
52     {
53         scratch_size = array_size / 2;
54         results = new double[scratch_size]();
55     }
56
57     // Gets wall time and start parallel block.
58     get_walltime(&start_time);
```



```

58 #pragma omp parallel
59 {
60     int threads, thread_id, start, end;
61     threads = omp_get_num_threads();
62     thread_id = omp_get_thread_num();
63     #pragma omp single
64     {
65         /** Sets variables shared between threads and handles the final
66         element in arrays first dependent on
67         * whether or not the scratch array size is even or odd. */
68         if (array_size % 2 == 1)
69         {
70             results[scratch_size - 1] = array[array_size - 1]; // Dynamic
71             Allocation
72             } else
73             {
74                 results[scratch_size - 1] = array[array_size - 1] + array[
75                 array_size - 2]; // Dynamic Allocation
76             }
77         // divide and conquer b-tree adds based on number of threads
78         num_elements = (int) ceil((double) (scratch_size) / (double) threads)
79         ;
80     }
81     // Compute start and end values for each thread
82     start = num_elements * thread_id;
83     end = start + num_elements;
84
85     // Ends at scratch array - 1 because final element is handled above.
86     for (start; start < end && start < scratch_size - 1; start++)
87     {
88         // start*2 is used to access the correct array elements:
89         // results[i] = array[2*i] + array[2*i + 1]
90         results[start] = array[start * 2] + array[start * 2 + 1];
91     }
92 }
93 // RECURSE! Scratch time is a dummy variable created in place of overloading
94 // the function call
95 answer = binary_sum(results, scratch_size, &scratch_time);
96 get_walltime(&end_time);
97 *time = end_time - start_time;
98 free(results); // frees dynamically allocated memory
99 return answer;
100 }
101
102 /** parallel_sum
103 *
104 * Calculates the sum of an array by dividing array into subsections and then
105 * summing those on different threads.
106 *
107 * @param *array - the array to be computed
108 * @param array_size - necessary for proper assignment
109 * @param *time - reference to a variable to store the walltime
110 *
111 * @return - sum of array as a double
112 */
113 double parallel_sum(double *array, int array_size, long double *time)
114 {
115     double *scratch, sum;

```

```

112     long double start_time = 0.0, end_time = 0.0;
113     int scratch_size, num_elements;
114     sum = 0.0;
115     get_walltime(&start_time);
116     #pragma omp parallel
117     {
118         int threads, thread_id, start, end;
119         threads = omp_get_num_threads();
120         thread_id = omp_get_thread_num();
121         #pragma omp single
122         {
123             /** Allocate variables necessary to share between threads. Scratch
124             space is allocated based on the
125             * number of threads obtained by OpenMP (one array element per thread
126             ). */
127             scratch = new double[threads](); // Dynamic allocation
128             scratch_size = threads;
129
130             // Divide array up into a set number of elements per thread.
131             num_elements = (int) ceil((double) array_size / (double) threads);
132
133             // Set start and endpoints for each thread
134             start = num_elements * thread_id;
135             end = start + num_elements;
136
137             // Sum up elements into scratch array
138             for (start; start < end && start < array_size; start++)
139             {
140                 scratch[thread_id] += array[start];
141             }
142
143             // Go through scratch array and calculate final sum.
144             for (int i = 0; i < scratch_size; i++)
145             {
146                 sum += scratch[i];
147             }
148             free(scratch); // Free dynamically allocated memory
149             get_walltime(&end_time);
150             *time = end_time - start_time;
151             return sum;
152         }
153     }
154
155 /** fast_sum
156 *
157 * Calculates the sum of an array using the OpenMP reduction method.
158 *
159 * @param *array - the array to be computed
160 * @param array_size - necessary for proper assignment
161 * @param *time - reference to a variable to store the walltime
162 *
163 * @return - sum of array as a double
164 */
165 double fast_sum(double *array, int array_size, long double *time)
166 {
167     double sum;
168     long double start_time = 0.0, end_time = 0.0;
169     sum = 0.0;

```

```

170     get_walltime(&start_time);
171     #pragma omp parallel for reduction(+:sum)
172     // declares sum to be the accumulator variable for all threads
173     for (int i = 0; i < array_size; i++)
174     {
175         sum += array[i];
176     }
177     get_walltime(&end_time);
178     *time = end_time - start_time;
179     return sum;
180 }
181
182
183 int main()
184 {
185     int set_threads = 10; // Threads to be requested by OpenMP
186     char filename[11];
187     sprintf(filename, "log-%02d.txt", set_threads);
188     omp_set_num_threads(set_threads); // Set global variable (shell script did
189     // not appear to actually set value)
190
191     std::ofstream log;
192     log.open(filename, std::ofstream::out | std::ofstream::app); // Open file
193     log << "ArraySize\tBinarySumTime\tBinarySumTotal\tParallelSumTime\t
194     tParallelSumTotal\tFastSumTime\t"
195     "FastSumTotal\tBinParDiff\tBinFastDiff\tParFastDiff\n";
196
197     for (int array_size = 1000000; array_size <= 10000000; array_size += 100000)
198     {
199         auto *x = new double[array_size](); // Dynamic allocation
200         long double parallel_time = 0.0, fast_time = 0.0, binary_time = 0.0;
201         for (int i = 0; i < array_size; i++)
202         {
203             x[i] = (double) (rand()) / (double) (RANDMAX) * 5.0;
204         }
205         double binary_total = binary_sum(x, array_size, &binary_time);
206         double parallel_total = parallel_sum(x, array_size, &parallel_time);
207         double fast_total = fast_sum(x, array_size, &fast_time);
208         log << array_size << "\t" << binary_time << "\t" << binary_total << "\t"
209         << parallel_time << "\t"
210         << parallel_total << "\t" << fast_time << "\t" << fast_total << "\t"
211         << binary_total - parallel_total
212         << "\t" << binary_total - fast_total << "\t" << parallel_total -
213         fast_total << std::endl;
214         free(x); // Free dynamically allocated memory
215     }
216     log.close(); // Close file
217     return 0;
218 }

```