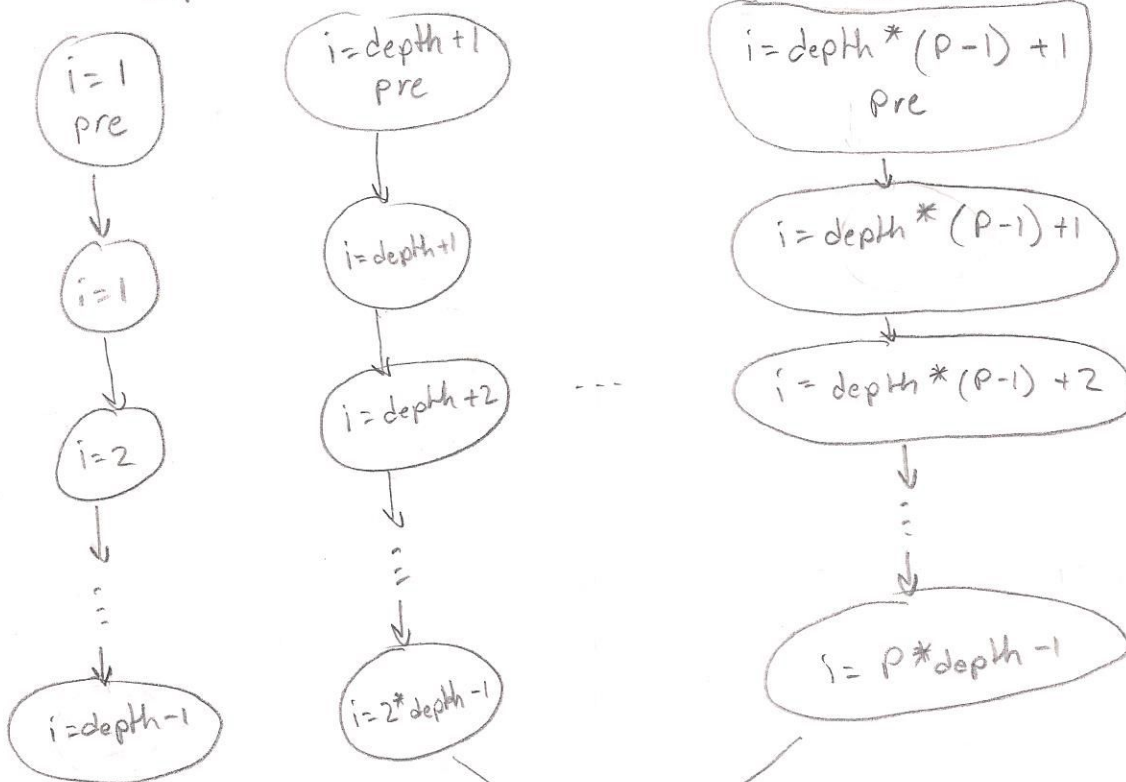


All tasks are dependent on variable "result"

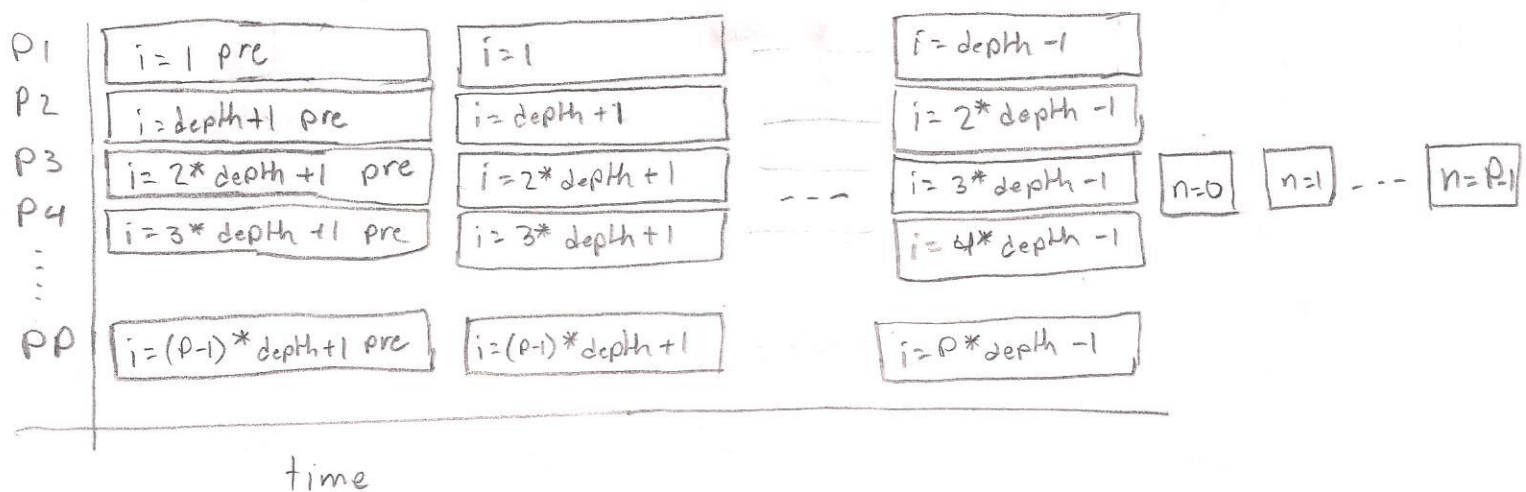
width = 1
CP = n
work = n

depth = floor(n/P) ← this is the number of additions each processor will do



each thread is in charge of "depth" additions then one processor adds up the partial sums

width = P
CP = depth + P
work = n + P



Each thread is in charge of "depth" number of sums. Then one thread (P3) is in charge of adding up the partial sums. The pre task for each thread is the initial set of the result variable

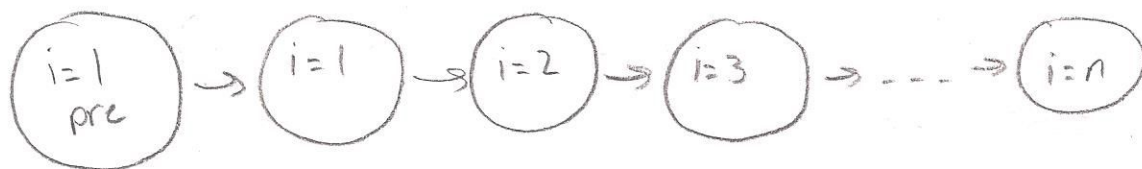
3.2) `reduce<int, max>` - yes. It would find the max element in each section of the array then compare the max for each section against each other to find the whole max

`reduce<std::string, concat>` - yes, He also would concat each section then concat each section together, it would maintain proper order

`reduce<float, sum>` - yes, code is templeted, and int is similar to float & would work the same way

`reduce<float, max>` - yes, similar to `reduce<int, max>`

4)



all tasks are dependent on the $i-1$ iteration. This is because to calculate $pr[i]$, you need $pr[i-1]$

width = 1
CP = $n+1$
work = $n+1$

Parallel version

Similar to the reduce function, the array will be split into smaller chunks and prefix sum will be computed then the last sum of each chunk is added into the next chunk

Pass 1	ar	pr	ar	pr	ar	pr
	0	0	6	0	12	0
	1	0	7	6	13	12
	2	1	8	13	14	25
	3	3	9	21	15	39
	4	6	10	30	16	54
	5	10	11	40	17	70
		15		51		87
	thread 1 offset = 15		thread 2 offset = 51		thread 3 offset = 87	

do partial prefix sum for each chunk of the data in this case 3 chunks. the offset is the largest partial sum of that chunk. This is step 1

Pass 2	pr	fin	pr	fin	pr	fin
	0	0	0	15	0	66
	0	0	6	21	12	78
	1	1	13	28	25	91
	3	3	21	36	39	105
	6	6	30	45	54	120
	10	10	40	55	70	136
	15	15	51	66	87	153
	error = 0		error = 15		error = 51 + 15 = 66	

Now add all the offsets from the previous chunks to each member of the partial sum the sum of the offsets is the error. This summation is independent and can be threaded similar to step 1. This is step 2 and final step

5) mergesort super highlevel pseudocode

mergesort(arr)

if $\text{len}(\text{arr}) \leq 1$:
return arr

left = left half of arr
right = right half of arr

left = mergesort(left)
right = mergesort(right)

return merge(left, right)

merge(left, right)

result = empty array

while left is not empty
and right is not empty:

if $\text{left}[0] \leq \text{right}[0]$:
result += left[0]
rem left[0]

else:
result += right[0]
rem right[0]

if left is not empty
result += left

if right is not empty
result += right

return result

an array of size 1 is sorted

split into 2 parts

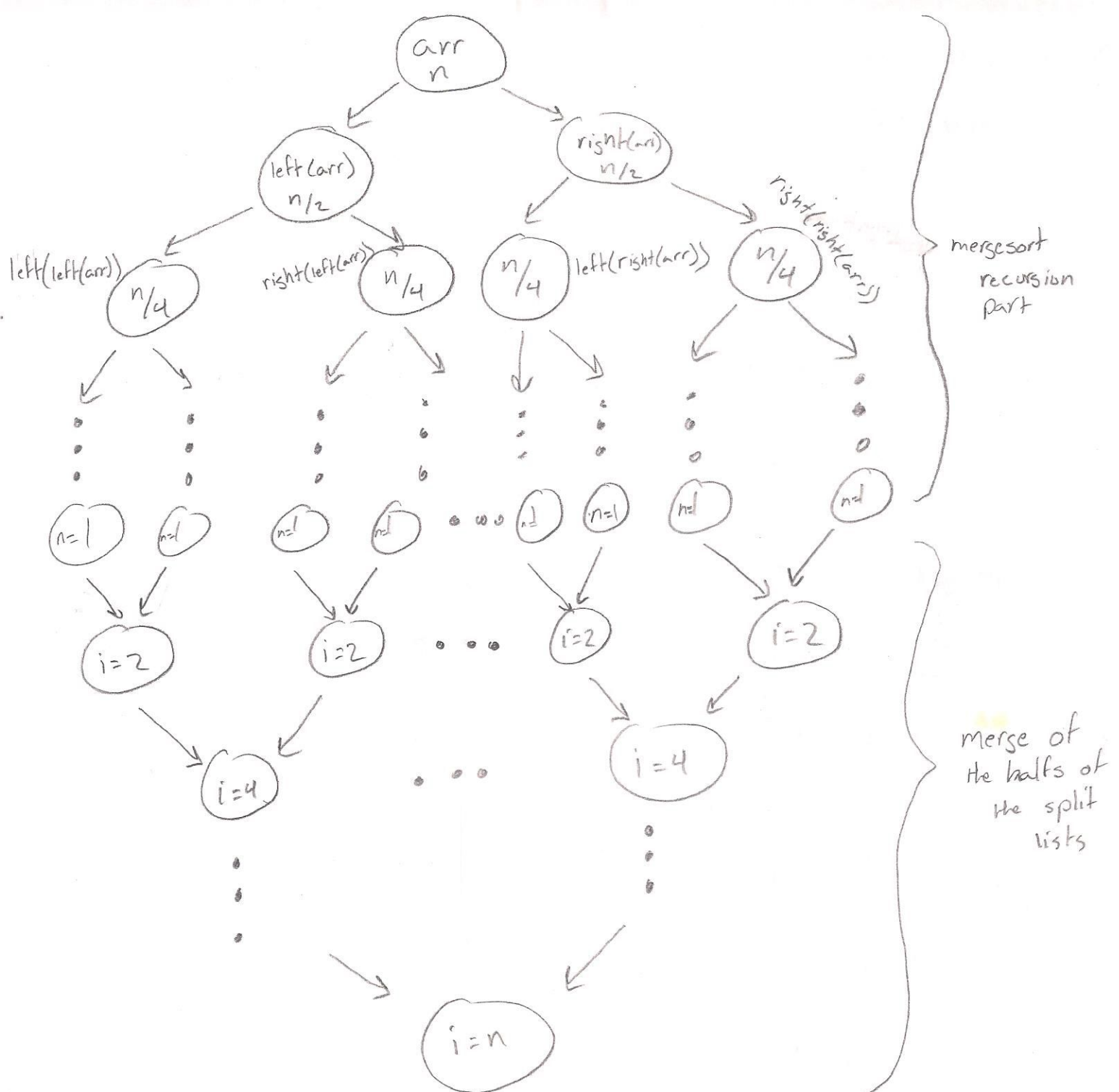
sort both left + right parts
recursively

merge the two parts in order

create empty result array

go through both left + right sides and
place the smaller of the first element into
the result array

place the rest of the elements from
either list into the result array

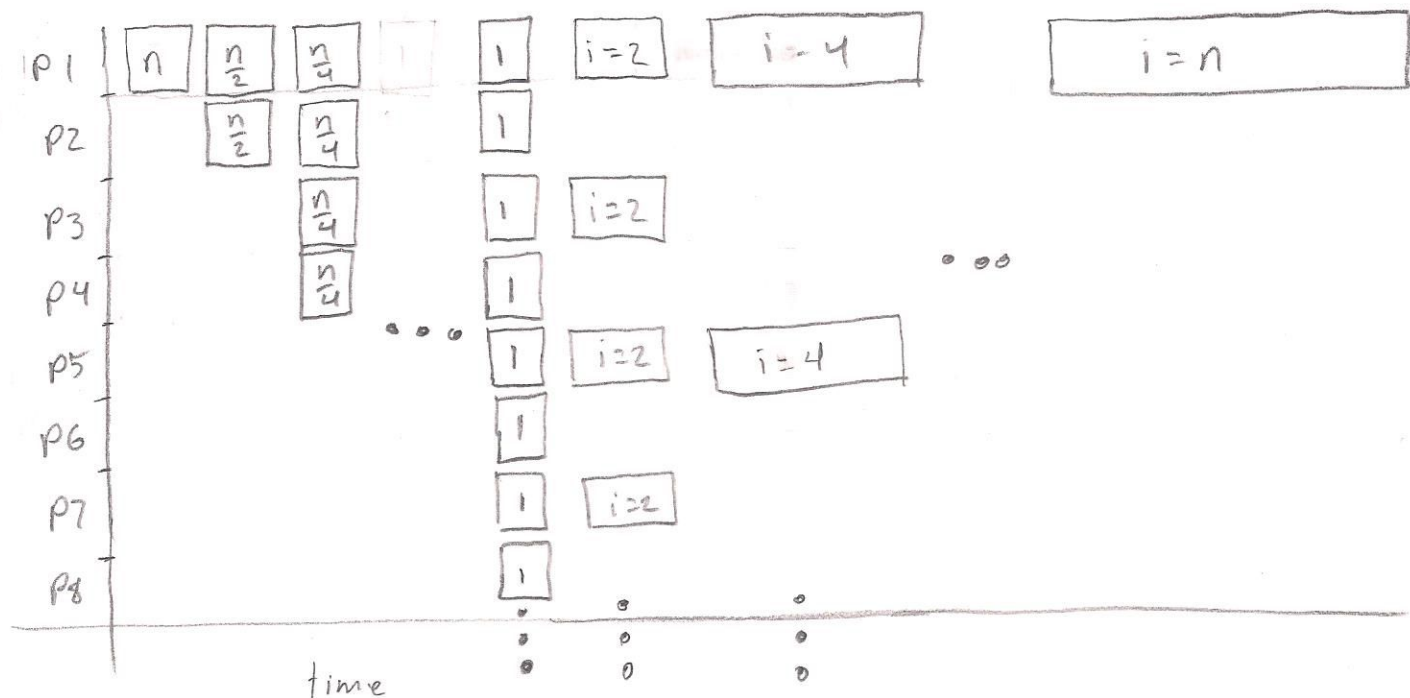


The first half of the dependency tree is the recursive part of merge sort() function. It splits the array into 2 subarrays. Then these subarrays are sorted and merged in the second half of the dependency graph. The merge part of the dependency graph takes more time to compute.

$$\text{work} = 3n - 2$$

$$\text{CP} = 2 \log_2(n) + 1$$

$$\text{width} = 1$$



Diamond shape? If we don't have any penalty for spawning threads, then each recursive call would spawn 2 threads: one for each half of its current thread. There is very little "work" happening other than the creation of the threads. That is why the bars here are thin. When the merger of two subarrays happens, one thread will handle it but because it has to go through 2 subarrays, the time taken to sort will double as the recursion stack collapses to the end condition ' $i=n$ '.

To make the code more parallel, I would remove some of the initial recursion so there is little dependency at the beginning of the run. For example with a P of 4, the first 2 "layers" of the Gantt chart / dependency graph will be gone. Then the processor mergesorts the 4 halves of the array. At the end a thread sorts the 2 halves that are left.

```
1
2 // PLEASE COMPILE WITH g++ reduce_par.cpp -lpthread
3
4 #include <memory>
5 #include <iostream>
6 #include <vector>
7 #include <thread>
8 #include <cassert>
9 #include <string>
10 // #include <numeric>
11 // #include <execution>
12 #include <parallel/numeric>
13
14 template <typename T>
15 T sum(T a, T b) {
16     return a + b;
17 }
18
19 template <typename T>
20 T max(T a, T b) {
21     return a > b ? a : b;
22 }
23
24 // template <typename T, typename = std::enable_if_t<std::is_same_v<T,
25 // std::string>>>
26 template <typename T, typename = std::enable_if_t<std::is_same<T,
27 std::string>::value>>
28 constexpr auto concat = sum<T>;
29
30 // Code from the assignment. The last parameter is needed
31 // because F is a type, even though it can be modified to only
32 // have 2 parameters, I believe this is easier on the head
33 template <typename T, typename F>
34 T reduce_sin(T* array, size_t n, F op) {
35     T result = array[0];
36     for (int i=1; i<n; ++i)
37         result = op(result, array[i]);
38     return result;
39 }
40
41 // I do not want to deal with promises or futures, so Im doing this.
42 /// Partially reduces the array from [start, end) and places it in
43 /// result.
44 /// @note The *result variable does not have a race condidtion as it
45 /// is only being written by one thread.
46 template <typename T, typename F>
47 void partial_reduce(T* start, T* end, T* result, F op) {
48     *result = reduce_sin(start, std::distance(start, end), op);
49 }
50
51 /// Parallel reduce
52 template <typename T, typename F>
53 T reduce_par (T* array, size_t n, size_t p, F op) {
54     // create nessicary structures
55     std::unique_ptr<T[]> data{ new T[p] };
56     std::vector<std::thread> threads{ };
57
58     // calculate parallelism
59     size_t depth = n / p;
60     T* start = array;
```

```

59     T* end = array + depth;
60
61     // the actual parallel code
62     for (size_t i = 0; i < p - 1; ++i) {
63         // The <T, F> is needed because partial_reduce is a templated function so
we need
64         // to thread the <T, F> version of the function
65         threads.emplace_back(partial_reduce<T, F>, start, end, &data[i], op);
66         start = end;
67         end += depth;
68     }
69
70     // the last one (this thread will do it so we dont spin on waiting for them
to join)
71     partial_reduce(start, array + n, &data[p - 1], op);
72
73     // wait for them to finish
74     for (auto& t : threads)
75         t.join();
76
77     // final sum
78     return reduce_sin(data.get(), p, op);
79 }
80
81
82 // options:
83 // 1: reduce<int, sum>
84 // 2: reduce<int, max>
85 // 3: reduce<std::string, concat>
86 // 4: reduce<float, sum>
87 // 5: reduce<float, max>
88 #define OPTION 1
89 #define MAX 500'000'007
90 // @warning, the only one that doesnt really work very well is option 4.
Adding up
91 // more than 1mil values gives you two different answers and I'm mostly
certain its
92 // because of FP arithmetic
93
94
95 #if OPTION == 1
96     #define NUMERIC
97     using NUM_TYPE = int;
98     template <typename T> constexpr auto op_t = sum<T>;
99 #elif OPTION == 2
100     #define NUMERIC
101     using NUM_TYPE = int;
102     template <typename T> constexpr auto op_t = max<T>;
103 #elif OPTION == 3
104     template <typename T> constexpr auto op_t = concat<T>;
105 #elif OPTION == 4
106     #define NUMERIC
107     using NUM_TYPE = float;
108     template <typename T> constexpr auto op_t = sum<T>;
109 #elif OPTION == 5
110     #define NUMERIC
111     using NUM_TYPE = float;
112     template <typename T> constexpr auto op_t = max<T>;
113 #else
114     #error "[E] Valid option not chosen"

```



```

115 #endif
116
117
118 int main() {
119
120 #ifdef NUMERIC
121
122     constexpr unsigned long long LG_NUM = MAX;
123     std::vector<NUM_TYPE> data{};
124     data.push_back(3); // this is a "valid" dummy data for max
125
126     for (unsigned long long i = 1; i < LG_NUM + 1; ++i) {
127         data.push_back(1); // sums of ones
128     }
129
130 #else
131
132     char temp[2] = { '\0', '\0' };
133     std::vector<std::string> data{};
134     for (int j = 0; j < 3; ++j) {
135         for (unsigned long long i = 0; i < 26; ++i) {
136             // 97-122
137             temp[0] = i + 97;
138             data.emplace_back(temp);
139         }
140     }
141
142 #endif
143
144     auto op = op_t<decltype(data)::value_type>;
145
146     auto begin1 = std::chrono::high_resolution_clock::now(); //// AHHHHHHH
147     auto accum1 = reduce_sin(data.data(), data.size(), op);
148     auto end1 = std::chrono::high_resolution_clock::now(); //// AHHHHHHH
149
150     std::chrono::duration<double, std::milli> elapse1{ end1 - begin1 };
151     std::cout << "Serial: " << elapse1.count() << std::endl;
152
153     auto begin2 = std::chrono::high_resolution_clock::now(); //// AHHHHHHH
154     auto accum2 = reduce_par(data.data(), data.size(), 8, op);
155     auto end2 = std::chrono::high_resolution_clock::now(); //// AHHHHHHH
156
157     std::chrono::duration<double, std::milli> elapse2{ end2 - begin2 };
158     std::cout << "Parallel: " << elapse2.count() << std::endl;
159
160     // Turns out that neither gnu g++ or clang++ supports the parallelism TS,
kms
161     // https://youtu.be/McJrc2uxbKI?t=571 MSVC ahead of the game *slow claps*
162     // nvm: https://godbolt.org/z/TCGaze The internet video lies
163     // auto begin3 = std::chrono::high_resolution_clock::now();
164     // double result = std::reduce(std::execution::par, sum.begin(),
sum.end());
165     // auto end3 = std::chrono::high_resolution_clock::now();
166
167     // std::chrono::duration<double, std::milli> elapse3{ end3 - begin3 };
168     // std::cout << "STL: " << elapse3.count();
169
170     // hehe https://gcc.gnu.org/onlinedocs/libstdc++/manual/parallel_mode.html
171     // after seeing the results, another lie. Why does the docs lie?????
172     // Serial: 1512.65

```

```
173 // Parallel: 351.756
174 // GNU: 1562.03
175 // The sum is: 500000010
176 // #if OPTION == 1
177 // auto begin4 = std::chrono::high_resolution_clock::now();
178 // double accum4 = __gnu_parallel::accumulate(data.begin(), data.end(), 0);
179 // auto end4 = std::chrono::high_resolution_clock::now();
180
181 // std::chrono::duration<double, std::milli> elapse4{ end4 - begin4 };
182 // std::cout << "GNU: " << elapse4.count() << std::endl;
183 // #endif
184
185 if (accum1 == accum2)
186     std::cout << "The result is: " << accum1 << std::endl;
187 else
188     std::cout << "The result were different, you screwed up..." << std::endl
<< accum1 << " " << accum2 << std::endl;
189 }
```