# Question 3

Using the C Programming language, write a program that sums an array of 50 elements. Next, optimize the code using loop unrolling. Loop unrolling is a program transformation that reduces the number of iterations for a loop by increasing the number of elements computed on each iteration. Generate a graph of performance improvement. Tip: Figure 5.17 in the textbook provides an example of a graph depicting performance improvements associated with loop unrolling. [30 marks]

Original program sumarray:

```
 8    int sumarray(int input[AMOUNT])
 9    {
10        int i;
11        int acc = 0;
12        for (i = 0; i < AMOUNT; i++)
13        {
14            int g;
15            for (g = 0; g < DELAY; g++);
16            acc += input[i];
17        }
18        return acc;
19    }
```

sumarray1 which uses 2 x 1 unrolling:

```
21    int sumarray1(int input[AMOUNT])
22    {
23        int i;
24        int acc = 0;
25        int limit = AMOUNT-1;
26
27        for (i = 0; i < limit; i+=2)
28        {
29            int g;
30            for (g = 0; g < DELAY; g++);
31            acc = (acc + input[i]) + input[i+1];
32        }
33
34        for (; i < AMOUNT; i++)
35            acc += input[i];
36
37        return acc;
38    }
```

sumarray2 which uses 2 x 1 unrolling with re-association:

```
40    int sumarray2(int input[AMOUNT])
41    {
42        int i;
43        int acc = 0;
44        int limit = AMOUNT-1;
45
46        for (i = 0; i < limit; i+=2)
47        {
48            int g;
49            for (g = 0; g < DELAY; g++);
50            acc += (input[i] + input[i+1]);
51        }
52
53        for (; i < AMOUNT; i++)
54            acc += input[i];
55
56        return acc;
57    }
```

sumarray3 which uses 5 x 1 unrolling and re-association:

```
59    int sumarray3(int input[AMOUNT])
60    {
61        int i;
62        int acc = 0;
63        int limit = AMOUNT-4;
64
65        for (i = 0; i < limit; i+=5)
66        {
67            int g;
68            for (g = 0; g < DELAY; g++);
69            acc += (input[i] + input[i+1]);
70            acc += (input[i+2] + input[i+3]);
71            acc += input[i+4];
72        }
73
74        for (; i < AMOUNT; i++)
75            acc += input[i];
76
77        return acc;
78    }
```

sumarray4 which uses 5 x 5 unrolling:

```c
int sumarray4(int input[AMOUNT])
{
    int i;
    int acc0 = 0;
    int acc1 = 0;
    int acc2 = 0;
    int acc3 = 0;
    int acc4 = 0;
    int limit = AMOUNT-4;

    for (i = 0; i < limit; i+=5)
    {
        int g;
        for (g = 0; g < DELAY; g++);
        acc0 += input[i];
        acc1 += input[i+1];
        acc2 += input[i+2];
        acc3 += input[i+3];
        acc4 += input[i+4];
    }

    for (; i < AMOUNT; i++)
        acc0 += input[i];

    return acc0 + acc1 + acc2 + acc3 + acc4;
}
```

sumarray5 which uses 10 x 10 unrolling:

```
107    int sumarray5(int input[AMOUNT])
108    {
109        int i;
110        int acc0 = 0;
111        int acc1 = 0;
112        int acc2 = 0;
113        int acc3 = 0;
114        int acc4 = 0;
115        int acc5 = 0;
116        int acc6 = 0;
117        int acc7 = 0;
118        int acc8 = 0;
119        int acc9 = 0;
120        int limit = AMOUNT-4;
121
122        for (i = 0; i < limit; i+=10)
123        {
124            int g;
125            for (g = 0; g < DELAY; g++);
126            acc0 += input[i];
127            acc1 += input[i+1];
128            acc2 += input[i+2];
129            acc3 += input[i+3];
130            acc4 += input[i+4];
131            acc5 += input[i+5];
132            acc6 += input[i+6];
133            acc7 += input[i+7];
134            acc8 += input[i+8];
135            acc9 += input[i+9];
136        }
137
138        for (; i < AMOUNT; i++)
139            acc0 += input[i];
140
141        return (acc0 + acc1 + acc2 + acc3 + acc4 + acc5 + acc6 + acc7 + acc8 + acc9);
142    }
```

main:

```
1    #include <stdio.h>
2    #include <stdlib.h>
3    #include <time.h>
4    #define MAX 1000
5    #define AMOUNT 50
6    #define DELAY 0xfffffff
144  int main ()
145  {
146      int array[AMOUNT];
147      int i;
148      srand(time(NULL));
149      for (i = 0; i < AMOUNT; i++)
150          array[i] = rand() % MAX; // use remainder operator to limit the size.
151
152      int sum = sumarray(array);
153      printf("The sum of the values of the array is %d\n", sum);
154
155      int sum1 = sumarray1(array);
156      printf("The sum of the values of the array is %d\n", sum1);
157
158      int sum2 = sumarray2(array);
159      printf("The sum of the values of the array is %d\n", sum2);
160
161      int sum3 = sumarray3(array);
162      printf("The sum of the values of the array is %d\n", sum3);
163
164      int sum4 = sumarray4(array);
165      printf("The sum of the values of the array is %d\n", sum4);
166
167      int sum5 = sumarray5(array);
168      printf("The sum of the values of the array is %d\n", sum5);
169
170      exit (0);
171  }
```
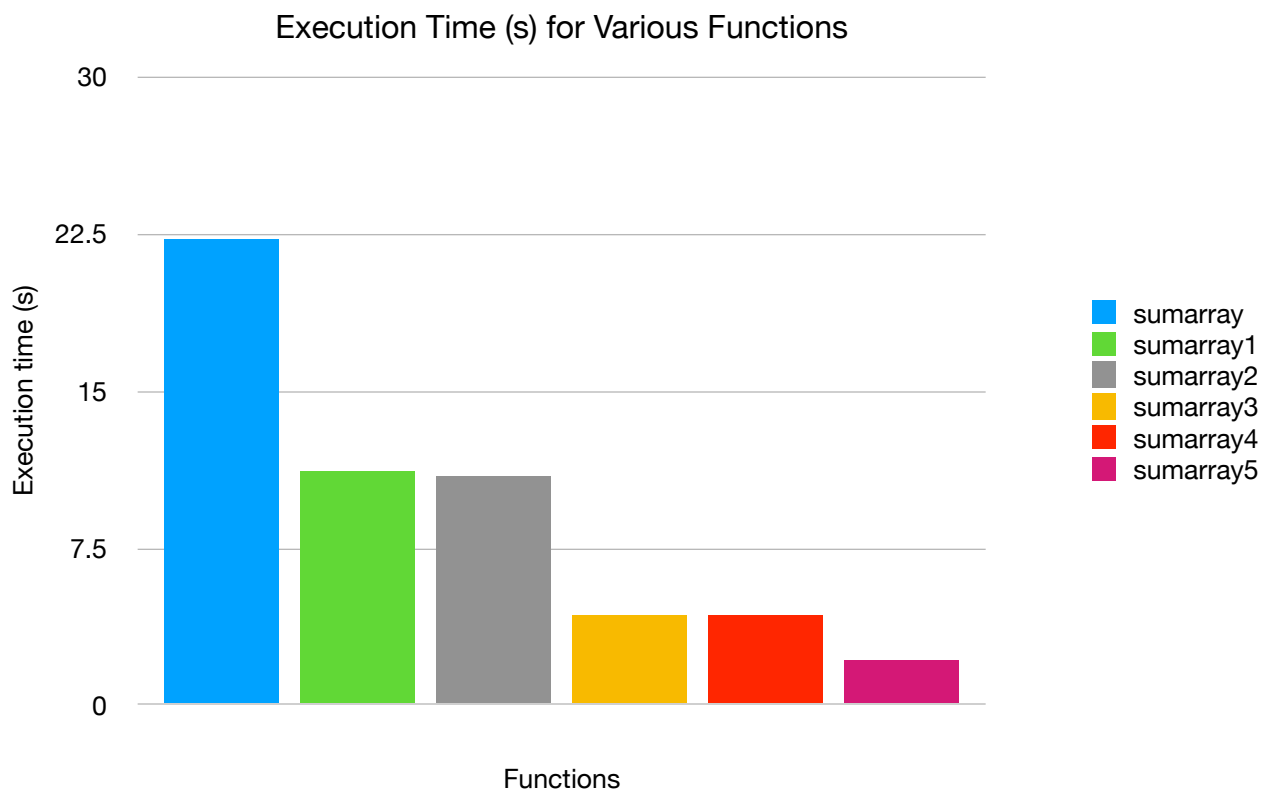
Execution results:

```
MichaelKuby@Michaels-iMac Question 3 % gcc-11 -O0 sumarray.c -o sumarray
MichaelKuby@Michaels-iMac Question 3 % ./sumarray
The sum of the values of the array is 24332
The sum of the values of the array is 24332
The sum of the values of the array is 24332
The sum of the values of the array is 24332
The sum of the values of the array is 24332
The sum of the values of the array is 24332
MichaelKuby@Michaels-iMac Question 3 % []
```

# Analysis

By compiling with -pg and testing with GPROF we see some interesting results:

```
Each sample counts as 0.01 seconds.
  %   cumulative   self              self     total
 time   seconds   seconds    calls   s/call   s/call  name
 40.66     22.31     22.31        1    22.31    22.31  sumarray
 20.38     33.49     11.18        1    11.18    11.18  sumarray1
 19.84     44.38     10.89        1    10.89    10.89  sumarray2
  7.85     48.68      4.31        1     4.31     4.31  sumarray4
  7.77     52.95      4.27        1     4.27     4.27  sumarray3
  3.86     55.07      2.12        1     2.12     2.12  sumarray5
```



Execution Time (s) for Various Functions

The third column "self seconds" shows us the total run time in seconds for each function. Our original code sumarray takes over 22 seconds to execute. The fastest of our 5 functions is sumarray5, which utilizes 10 x 10 unrolling. The difference here, rounded to the hundredth, is 22.31 / 2.12 = 10.52, meaning we've achieved a speed up of over a factor of 10. Pretty impressive.

From sumarray to sumarray1, which utilizes 2 x 1 unrolling with suboptimal association, we see an immediately drastic improvement of about a factor of 2. Somewhat surprisingly we notice that the reassociation technique utilized by sumarray2 offers almost no improvement over sumarray1. Nevertheless, it does add some improvement, and should be utilized.

From sumarray2 to sumarray3 we see another drastic improvement. sumarray3 uses 5 x 1 unrolling and the proper associations. At this point we have hit the latency bound of the hardware.

Interestingly, sumarray4 has a very slight regression in performance in attempting to use multiple accumulators. It's hard for me to pinpoint exactly why this is happening. My first guess was that the regression in performance was because the numerous number of variables cannot all be held in registers, causing variables to be stored in memory. This seemed unlikely given the fact that there are only 5 accumulators and x86-64 hardware has 16 general purpose registers. The results of sumarray5 go to show that the issue was not spillage.

sumarray5 is our top performer, breaking through the latency limit and getting what I would assume is somewhere close to the throughput limit. sumarray5 uses 10 x 10 unrolling, so it's speed is achieved via the use of a large number of accumulators, which exploits the functional capabilities of the systems hardware.