

## CS-525 Homework 1

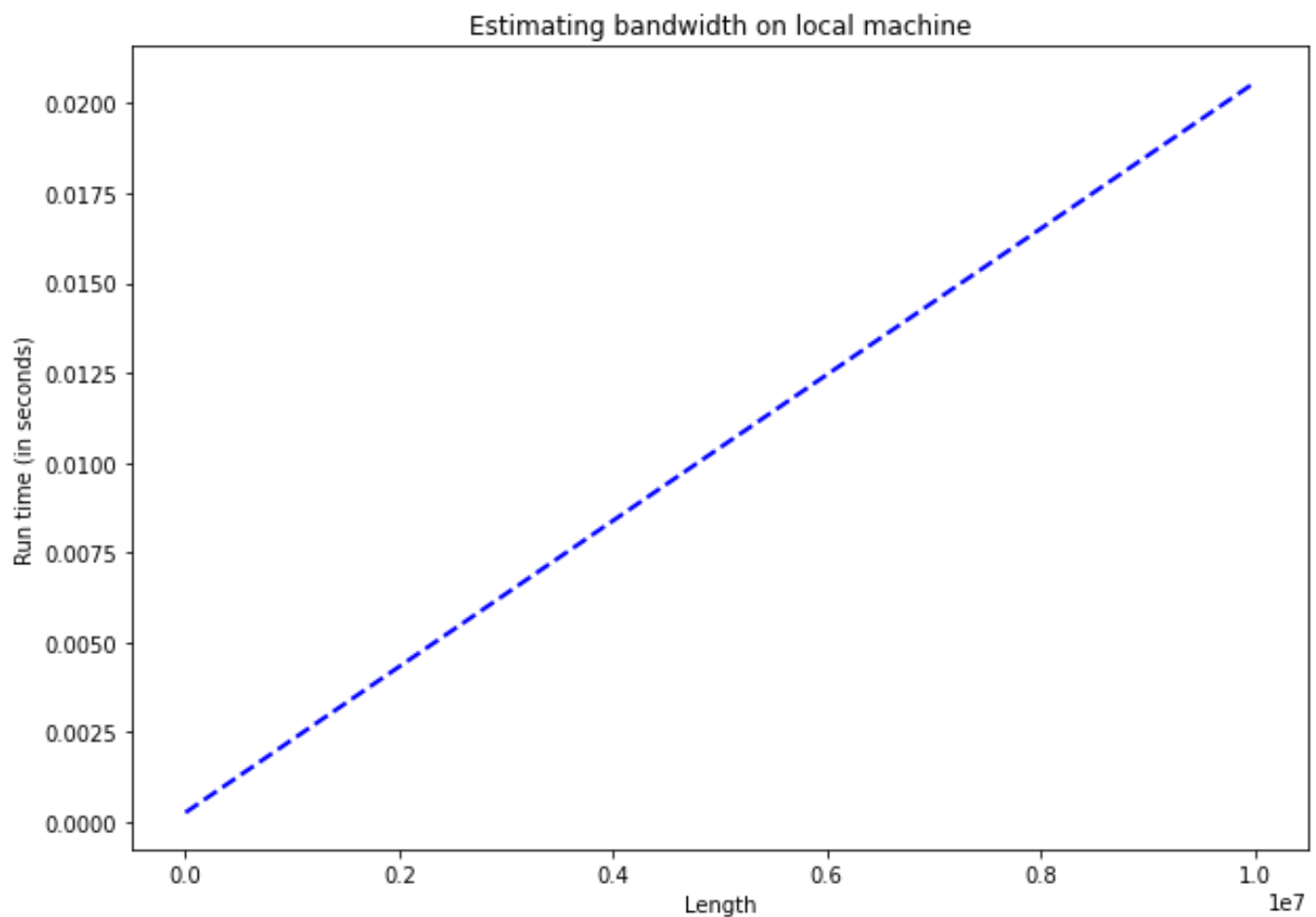
**Q1. Sol** - Following is the line chart of length vs runtime **recorded on local machine** for running the array on length varying from 1000 to 10000000. Here, the array is an integer array.

I calculated the slope of the line using python's NumPy function, polyfit, which gave the value as  $2.0313350303314353 \times 10^{-9}$ .

As Bandwidth is calculated as  $(\text{number of iterations} * \text{sizeof(int)} / \text{runtime})$ ,

The bandwidth of the local machine will be  $(1/2.0313350303314353 \times 10^{-9}) * 4 \text{ Bytes}$ , which is approximately 1.96 GB/sec.

Therefore, the **bandwidth of the local machine is 1.96 GB/sec**.



**Output of the code:**

Length: 1000 Time Measured:  $5 \times 10^{-6}$  seconds

Length: 10000 Time Measured:  $4.7 \times 10^{-5}$  seconds

Length: 100000 Time Measured: 0.000446 seconds

Length: 1000000 Time Measured: 0.002871 seconds

Length: 10000000 Time Measured: 0.02052 seconds

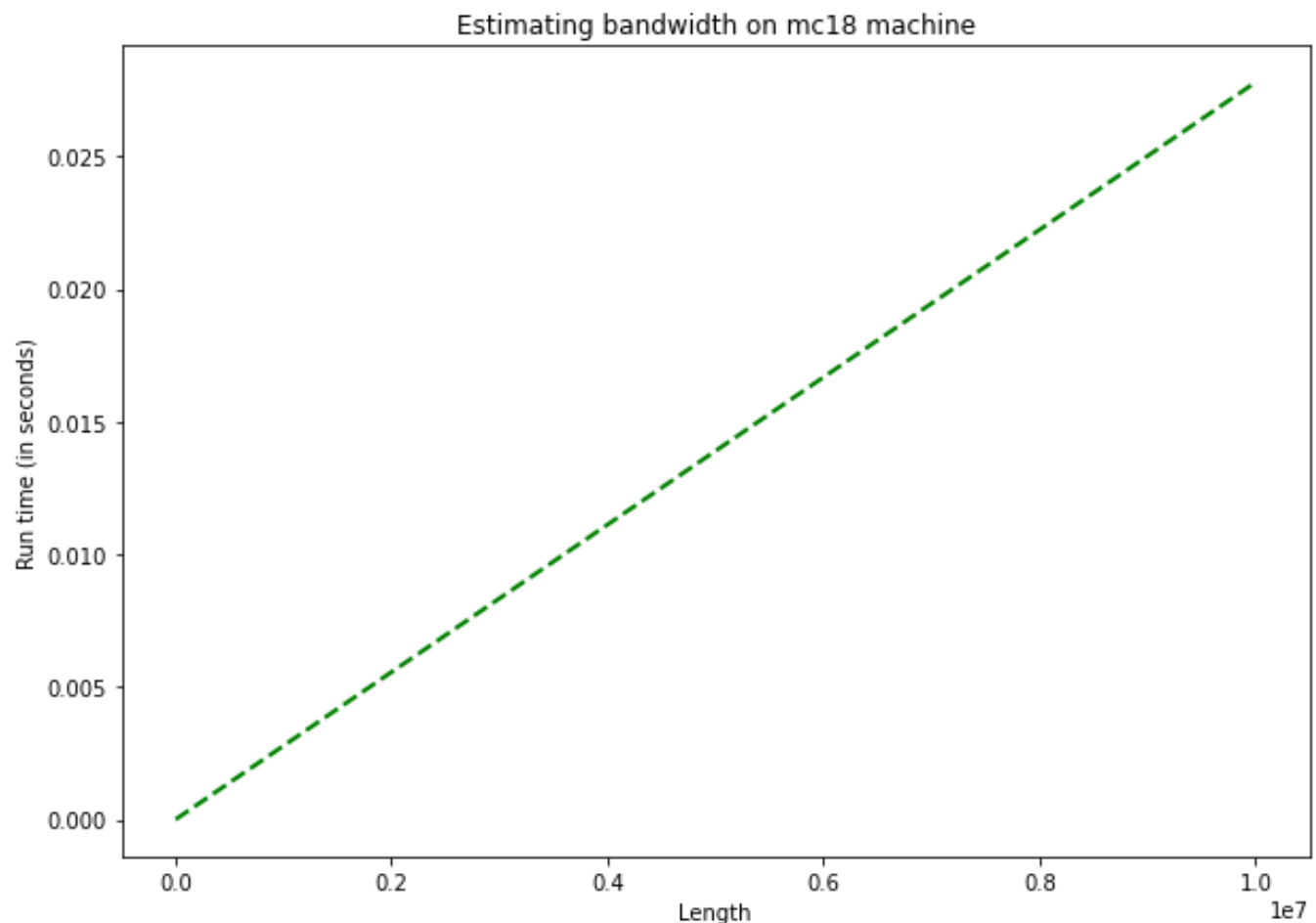
Following is the line chart of length vs runtime **recorded on mc18 machine** for running the array on length varying from 1000 to 10000000. Here, the array is an integer array.

I calculated the slope of the line using python's NumPy function, polyfit, which gave the value as  $2.779163334086171 \times 10^{-9}$ .

As Bandwidth is calculated as  $(\text{number of iterations} * \text{sizeof(int)} / \text{runtime})$ ,

The bandwidth of the local machine will be  $(1/2.779163334086171 \times 10^{-9}) * 4$  Bytes, which is approximately 1.4 GB/sec.

Therefore, the **bandwidth of the mc18 machine is 1.4 GB/sec.**



**Output of the code:**

Length: 1000 Time Measured: 4e-06 seconds

Length: 10000 Time Measured: 3.3e-05 seconds

Length: 100000 Time Measured: 0.000319 seconds

Length: 1000000 Time Measured: 0.003292 seconds

Length: 10000000 Time Measured: 0.027753 seconds

**Execute the code:**

Code will be found inside hw1/question1/ folder

**Compile:** g++ main.cpp -o ./run

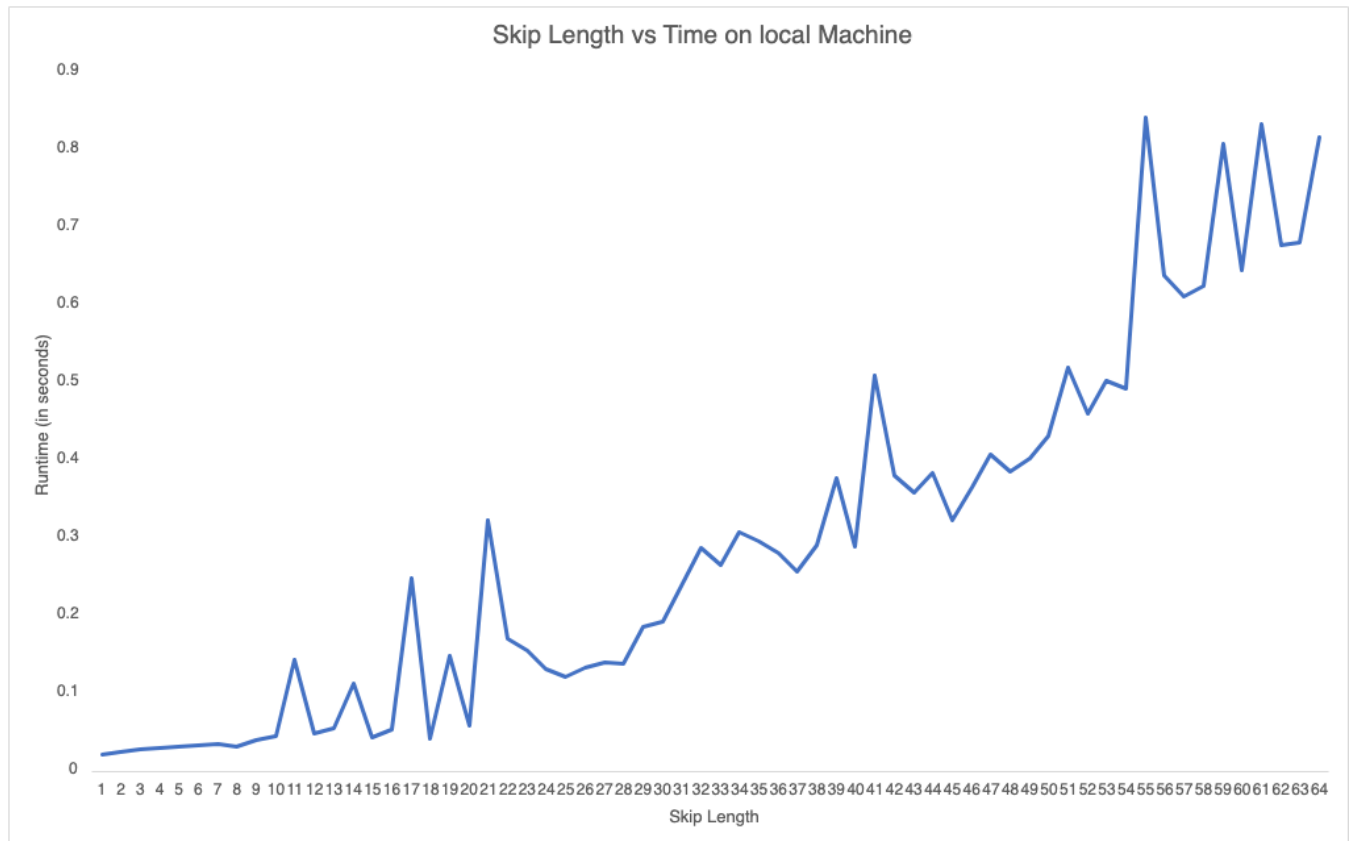
**Execute:** ./run

**Output:** output will be saved in output.txt file

**Q2. Sol** – Following is the line chart of skip length vs the runtime **recorded on local machine** for strided access on an array of length 10000000.

Here, we are trying to change the skip value from 1 to 64, and accordingly determining the time it takes for each execution.

We see multiple spikes/jumps in the graph. The First Spike occurs when the L1 cache size = skip size \* 4 Bytes (As integer array is used), hence leading to cache miss and increasing the runtime. Therefore, the cache line size is  $11 * 4 = 44$  Bytes as the first spike occurs on skip = 11.



**Output of the code:**

Skip	:	1	Time	Measured:	0.0215	second
Skip	:	2	Time	Measured:	0.0257	second
Skip	:	3	Time	Measured:	0.0287	second
Skip	:	4	Time	Measured:	0.0306	second
Skip	:	5	Time	Measured:	0.0311	second

Skip	:	6	Time	Measured:	0.0343	seconds
				24		
Skip	:	7	Time	Measured:	0.0352	seconds
				56		
Skip	:	8	Time	Measured:	0.0315	seconds
				55		
Skip	:	9	Time	Measured:	0.0399	seconds
				59		
Skip	:	10	Time	Measured:	0.0447	seconds
				96		
Skip	:	11	Time	Measured:	0.1433	seconds
				05		
Skip	:	12	Time	Measured:	0.0488	seconds
				47		
Skip	:	13	Time	Measured:	0.0563	seconds
				47		
Skip	:	14	Time	Measured:	0.1122	seconds
				99		
Skip	:	15	Time	Measured:	0.0432	seconds
				28		
Skip	:	16	Time	Measured:	0.0546	seconds
				85		
Skip	:	17	Time	Measured:	0.2481	seconds
				18		
Skip	:	18	Time	Measured:	0.0415	seconds
				28		
Skip	:	19	Time	Measured:	0.1479	seconds
				47		
Skip	:	20	Time	Measured:	0.0594	seconds
				52		
Skip	:	21	Time	Measured:	0.3225	seconds
				96		
Skip	:	22	Time	Measured:	0.1705	seconds
				89		
Skip	:	23	Time	Measured:	0.1553	seconds
				86		
Skip	:	24	Time	Measured:	0.1321	seconds
				85		
Skip	:	25	Time	Measured:	0.1214	seconds
				09		
Skip	:	26	Time	Measured:	0.1338	seconds
				9		
Skip	:	27	Time	Measured:	0.1408	seconds
				65		

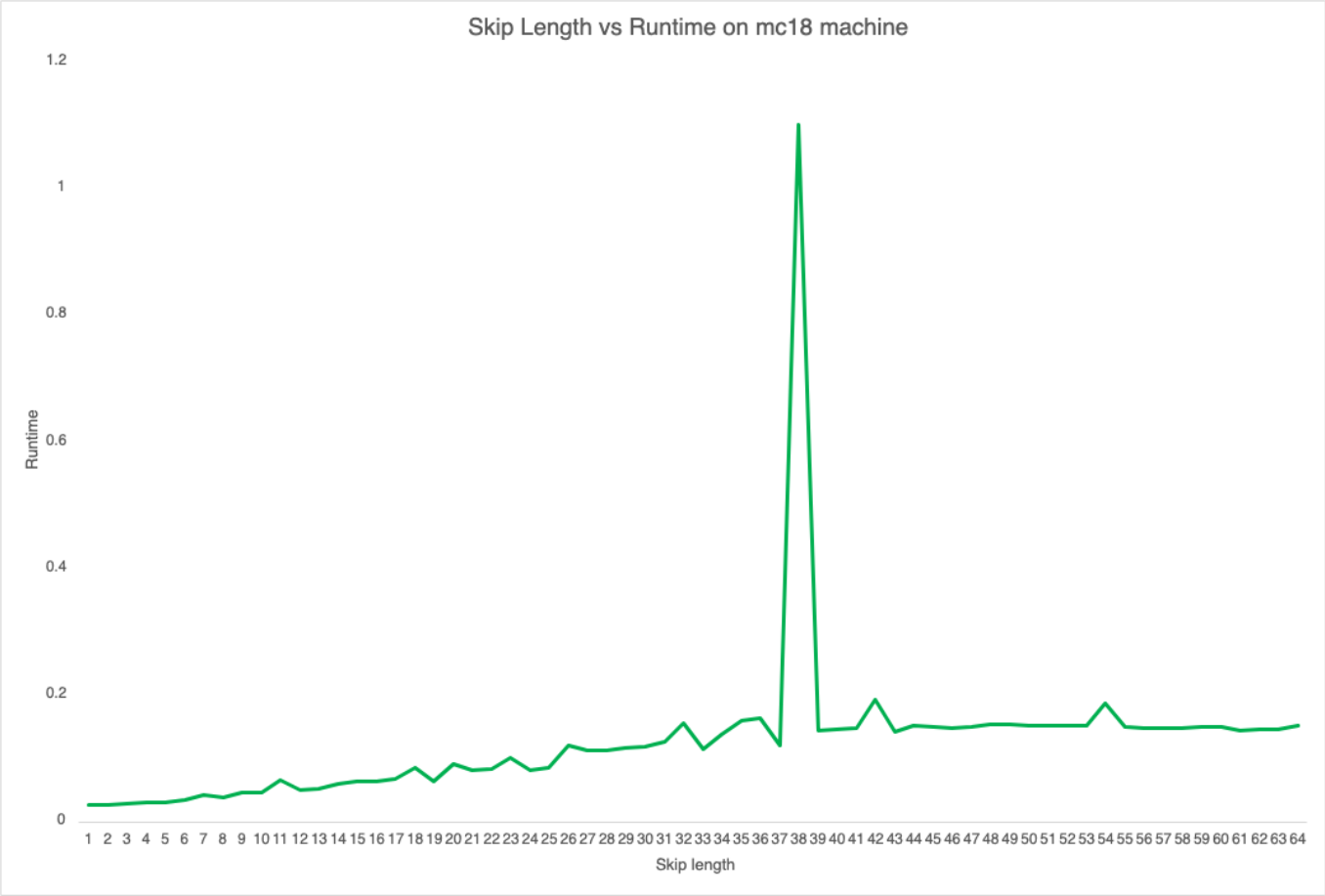
Skip	:	28	Time	Measured:	0.1388	seconds
Skip	:	29	Time	Measured:	0.1858	seconds
Skip	:	30	Time	Measured:	0.1929	seconds
Skip	:	31	Time	Measured:	0.2384	seconds
Skip	:	32	Time	Measured:	0.2871	seconds
Skip	:	33	Time	Measured:	0.2658	seconds
Skip	:	34	Time	Measured:	0.3083	seconds
Skip	:	35	Time	Measured:	0.2966	seconds
Skip	:	36	Time	Measured:	0.2800	seconds
Skip	:	37	Time	Measured:	0.2579	seconds
Skip	:	38	Time	Measured:	0.2917	seconds
Skip	:	39	Time	Measured:	0.3769	seconds
Skip	:	40	Time	Measured:	0.2895	seconds
Skip	:	41	Time	Measured:	0.5088	seconds
Skip	:	42	Time	Measured:	0.3807	seconds
Skip	:	43	Time	Measured:	0.3590	seconds
Skip	:	44	Time	Measured:	0.3842	seconds
Skip	:	45	Time	Measured:	0.3224	seconds
Skip	:	46	Time	Measured:	0.3647	seconds
Skip	:	47	Time	Measured:	0.4071	seconds
Skip	:	48	Time	Measured:	0.3863	seconds
Skip	:	49	Time	Measured:	0.4032	seconds

Skip	:	50	Time	Measured:	0.4316	second
					28	s
Skip	:	51	Time	Measured:	0.5198	second
					67	s
Skip	:	52	Time	Measured:	0.4606	second
					85	s
Skip	:	53	Time	Measured:	0.5029	second
					1	s
Skip	:	54	Time	Measured:	0.4920	second
					47	s
Skip	:	55	Time	Measured:	0.8408	second
					08	s
Skip	:	56	Time	Measured:	0.6386	second
					2	s
Skip	:	57	Time	Measured:	0.6104	second
					18	s
Skip	:	58	Time	Measured:	0.6254	second
					22	s
Skip	:	59	Time	Measured:	0.8071	second
					12	s
Skip	:	60	Time	Measured:	0.6445	second
					2	s
Skip	:	61	Time	Measured:	0.8328	second
					44	s
Skip	:	62	Time	Measured:	0.6778	second
					18	s
Skip	:	63	Time	Measured:	0.6803	second
					85	s
Skip	:	64	Time	Measured:	0.8167	second
					69	s

Following is the line chart of skip length vs the runtime **recorded on mc18 machine** for strided access on an array of length 10000000.

Here, we are trying to change the skip value from 1 to 64, and accordingly determining the time it takes for each execution.

We see multiple spikes/jumps in the graph. The First Spike occurs when the L1 cache size = skip size \* 4 Bytes (As integer array is used), hence leading to cache miss and increasing the runtime. Therefore, the cache line size is  $11 * 4 = 44$  Bytes as the first spike occurs on skip = 11.



Output of the code:

Skip	:	1	Time	Measur	0.0267	second
				ed:	48	s
Skip	:	2	Time	Measur	0.0277	second
				ed:	1	s
Skip	:	3	Time	Measur	0.0286	second
				ed:	36	s
Skip	:	4	Time	Measur	0.0304	second
				ed:	94	s
Skip	:	5	Time	Measur	0.0318	second
				ed:	65	s
Skip	:	6	Time	Measur	0.0344	second
				ed:	05	s
Skip	:	7	Time	Measur	0.0417	second
				ed:	22	s
Skip	:	8	Time	Measur	0.0386	second
				ed:	25	s
Skip	:	9	Time	Measur	0.0465	second
				ed:	53	s



Skip	:	10	Time	Measur	0.0469	second
				ed:	63	s
Skip	:	11	Time	Measur	0.0669	second
				ed:	8	s
Skip	:	12	Time	Measur	0.0499	second
				ed:	92	s
Skip	:	13	Time	Measur	0.0529	second
				ed:	23	s
Skip	:	14	Time	Measur	0.0610	second
				ed:	4	s
Skip	:	15	Time	Measur	0.0635	second
				ed:	65	s
Skip	:	16	Time	Measur	0.0639	second
				ed:	6	s
Skip	:	17	Time	Measur	0.0680	second
				ed:	81	s
Skip	:	18	Time	Measur	0.0855	second
				ed:	09	s
Skip	:	19	Time	Measur	0.0650	second
				ed:	79	s
Skip	:	20	Time	Measur	0.0910	second
				ed:	44	s
Skip	:	21	Time	Measur	0.0826	second
				ed:	85	s
Skip	:	22	Time	Measur	0.0831	second
				ed:	22	s
Skip	:	23	Time	Measur	0.1016	second
				ed:	76	s
Skip	:	24	Time	Measur	0.0823	second
				ed:	15	s
Skip	:	25	Time	Measur	0.0861	second
				ed:	49	s
Skip	:	26	Time	Measur	0.1215	second
				ed:	66	s
Skip	:	27	Time	Measur	0.1124	second
				ed:	89	s
Skip	:	28	Time	Measur	0.1127	second
				ed:	85	s
Skip	:	29	Time	Measur	0.1164	second
				ed:	53	s
Skip	:	30	Time	Measur	0.1183	second
				ed:	52	s
Skip	:	31	Time	Measur	0.1270	second
				ed:	22	s

Skip	:	32	Time	Measur	0.1569	second
				ed:	96	s
Skip	:	33	Time	Measur	0.1152	second
				ed:	73	s
Skip	:	34	Time	Measur	0.1379	second
				ed:	86	s
Skip	:	35	Time	Measur	0.1597	second
				ed:	17	s
Skip	:	36	Time	Measur	0.1640	second
				ed:	03	s
Skip	:	37	Time	Measur	0.1208	second
				ed:	83	s
Skip	:	38	Time	Measur	1.1012	second
				ed:	9	s
Skip	:	39	Time	Measur	0.1435	second
				ed:	22	s
Skip	:	40	Time	Measur	0.1457	second
				ed:	62	s
Skip	:	41	Time	Measur	0.1491	second
				ed:	05	s
Skip	:	42	Time	Measur	0.1942	second
				ed:	29	s
Skip	:	43	Time	Measur	0.1429	second
				ed:	55	s
Skip	:	44	Time	Measur		second
				ed:	0.1519	s
Skip	:	45	Time	Measur	0.1496	second
				ed:	32	s
Skip	:	46	Time	Measur	0.1483	second
				ed:	15	s
Skip	:	47	Time	Measur	0.1510	second
				ed:	74	s
Skip	:	48	Time	Measur	0.1532	second
				ed:	48	s
Skip	:	49	Time	Measur	0.1532	second
				ed:	97	s
Skip	:	50	Time	Measur	0.1526	second
				ed:	47	s
Skip	:	51	Time	Measur	0.1523	second
				ed:	86	s
Skip	:	52	Time	Measur	0.1523	second
				ed:	08	s
Skip	:	53	Time	Measur	0.1531	second
				ed:	21	s

Skip	:	54	Time	Measured:	0.1879	seconds
Skip	:	55	Time	Measured:	0.1505	seconds
Skip	:	56	Time	Measured:	0.1488	seconds
Skip	:	57	Time	Measured:	0.1490	seconds
Skip	:	58	Time	Measured:	0.1487	seconds
Skip	:	59	Time	Measured:	0.1499	seconds
Skip	:	60	Time	Measured:	0.1500	seconds
Skip	:	61	Time	Measured:	0.1453	seconds
Skip	:	62	Time	Measured:	0.1464	seconds
Skip	:	63	Time	Measured:	0.1461	seconds
Skip	:	64	Time	Measured:	0.1518	seconds

**Execute the code:**

Code will be found inside hw1/question2/ folder

**Compile:** `g++ main.cpp -o ./run`

**Execute:** `./run`

**Output:** output will be saved in output.txt file

### Q3. Sol -

Following is the output of loop unrolling from 1, 2, 3, 4, 5, 6 on **local machine**

Unroll count: 1 Time Measured: 0.015959 seconds

Unroll count: 2 Time Measured: 0.008089 seconds

Unroll count: 3 Time Measured: 0.007784 seconds

Unroll count: 4 Time Measured: 0.006991 seconds

Unroll count: 5 Time Measured: 0.006633 seconds

Unroll count: 6 Time Measured: 0.006363 seconds

Following is the output of unrolling from 1, 2, 3, 4, 5, 6 on **mc18 machine**

Unroll count: 1 Time Measured: 0.015094 seconds

Unroll count: 2 Time Measured: 0.013467 seconds

Unroll count: 3 Time Measured: 0.012783 seconds

Unroll count: 4 Time Measured: 0.012211 seconds

Unroll count: 5 Time Measured: 0.012026 seconds

Unroll count: 6 Time Measured: 0.011794 seconds

From the above results on both the platforms, we can comment that as the unrolling factor increases from 1 to 6, the time to execute to loop decreases. This is because:

- It reduces the loop overhead, i.e., incrementing and testing of the loop counter is reduced ,
- Improves the pipeline behavior by reducing jumps back to a loop's entry.
- It increases instruction-level parallelism, i.e., based on the unrolling factor, more instructions can be executed in parallel.

But, at the same time it increases the code size.

For example, consider the loop:

```
for (int i=0; i<length; i++) {  
    sum += arr[i];  
}
```

Here, let us start with `sum += arr[0]`. There will be a cache miss for `arr[0]` leading to accessing 4 words, .e., `arr[0]`, `arr[1]`, `arr[2]`, `arr[3]` from the memory. Now, all these 4 been in cache memory cannot be used in parallel, as the operation `sum += arr[i]` needs to be executed sequentially due to dependency on `sum` value.

Now, consider the loop unrolled twice:

```
for (int i=0; i<length; i=i+3) {  
    temp1 += arr[i];  
    temp2 += arr[i+1];  
    temp3 += arr[i+2];  
}  
int temp = temp1 + temp2 + temp3 ;
```

Here, let us start with `temp1 += arr[0]`. There will be a cache miss for `arr[0]` leading to accessing 4 words, i.e. `arr[0]`, `arr[1]`, `arr[2]`, `arr[3]` from the memory. Now, all these 4 been in cache memory can be used in parallel increasing instruction-level parallelism, as all the 4 instructions in the loop can be executed in parallel as there is no dependency between the variables. Also at each step, the counter is incremented by 3, which reduces the loop overhead, therefore, improving the performance.

#### **Execute the code:**

Code will be found inside `hw1/question3/` folder

**Compile:** `g++ main.cpp -o ./run`

**Execute:** `./run`

#### Q4. Sol -

Following is the output of impact of loop interchanging on **local machine**

Matrix Vector Multiply, **length :1000** Time Measured: 0.006192 seconds

Matrix Vector Multiply InterChange, length :1000 Time Measured: 0.0053 seconds

Matrix Vector Multiply, **length :10000** Time Measured: 0.356515 seconds

Matrix Vector Multiply InterChange, length :10000 Time Measured: 0.83891 seconds

For local machine, I was unable to perform Matrix Vector Multiplication for length of 100000 due to out of memory issue. Hence, I have not implemented it on local machine, but it works on mc18 machine.

Following is the output of impact of loop interchanging on **mc18 machine**

Matrix Vector Multiply, **length :1000** Time Measured: 0.006871 seconds

Matrix Vector Multiply InterChange, length :1000 Time Measured: 0.015166 seconds

Matrix Vector Multiply, **length :10000** Time Measured: 0.38019 seconds

Matrix Vector Multiply InterChange, length :10000 Time Measured: 1.65345 seconds

Matrix Vector Multiply, **length :100000** Time Measured: 41.8175 seconds

Matrix Vector Multiply InterChange, length :100000 Time Measured: 669.095 seconds

To explain the above scenario,

Consider case 1: Matrix Vector Multiply:

```
for (int i = 0; i < n; i++){  
    for (int j = 0; j < n; j++) {  
        result[i] += matrix[i][j]*b[j];  
    }  
}
```

Let is start with  $\text{result}[0] = \text{matrix}[0][0] * b[0]$ . This will have a cache miss at  $\text{matrix}[0][0]$  leading to accessing it in memory. As arrays are laid out as row major in C/C++, therefore, while accessing  $\text{matrix}[0][0]$ , we will get  $\text{matrix}[0][0], \text{matrix}[0][1], \text{matrix}[0][2], \text{matrix}[0][3]$ . Hence,

after the first cache miss, the next three iterations can be executed without cache miss. Same will be the case for vector b. There will be less cache miss leading to lesser runtime.

Now, consider case 2: Matric Vector Multiplication with loop interchange.

```
for (int j = 0; j < n; j++){  
  for (int i = 0; i < n; i++) {  
    result[i] += matrix[i][j]*b[j];  
  }  
}
```

Let is start with  $\text{result}[0] = \text{matrix}[0][0] * b[0]$ . Here too, we will get cache miss for  $\text{matrix}[0][0]$ , and we get  $\text{matrix}[0][0], \text{matrix}[0][1], \text{matrix}[0][2], \text{matrix}[0][3]$  words. But, however, in the second iteration, we will access  $\text{matrix}[1][0] * b[0]$ . This is a problem as there will be a miss for  $\text{matrix}[1][0]$  which will again give us 4 words,  $\text{matrix}[1][0], \text{matrix}[1][1], \text{matrix}[1][2], \text{matrix}[1][3]$ . Therefore, in every iteration, we get a cache miss as we are unable to use the bulk access, which leads to an increase in runtime as seen in the output.

Therefore, strided access is terrible, as we do not make use of bulk access and get worst performance.

#### **Execute the code:**

Code will be found inside hw1/question4/ folder

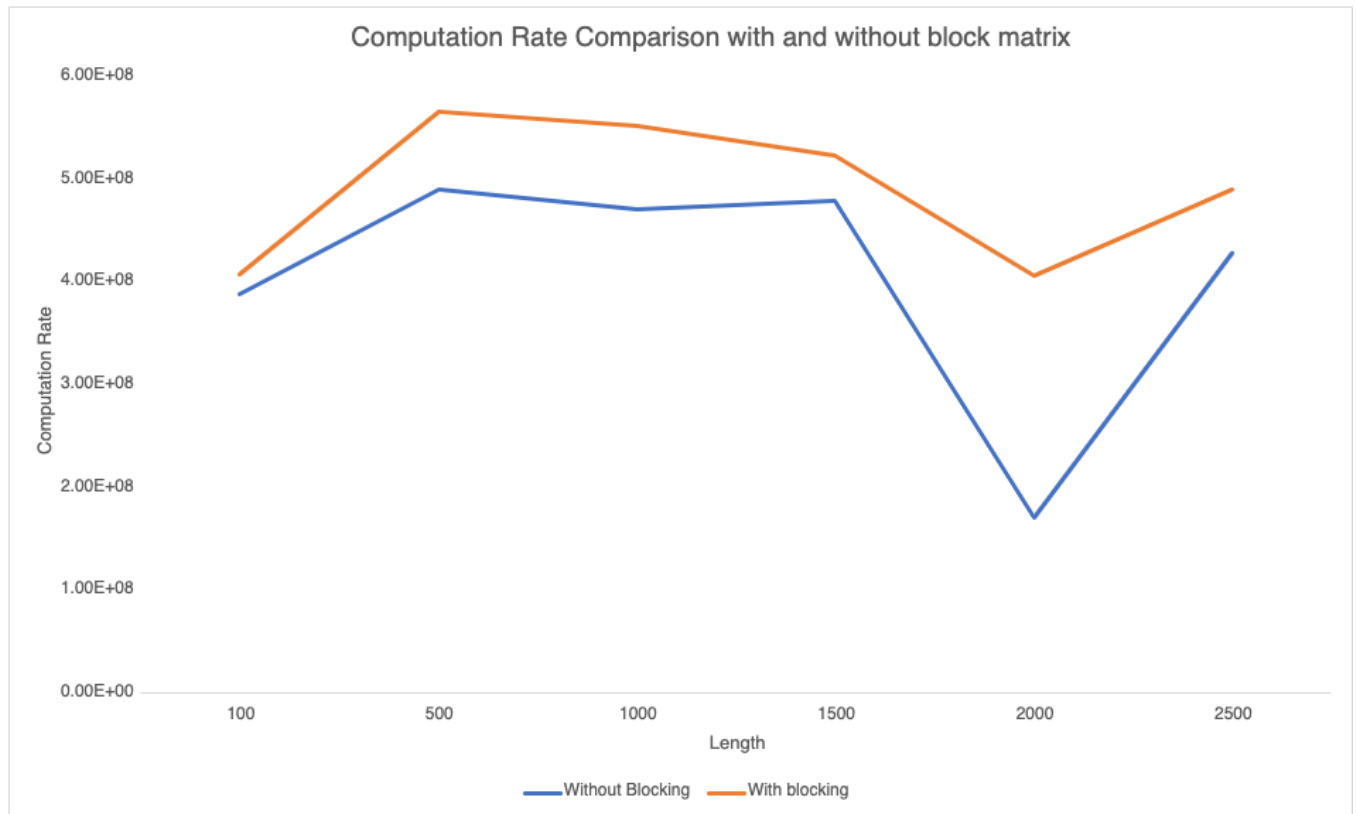
**Compile:** `g++ main.cpp -o ./run`

**Execute:** `./run`

**Output:** output will be saved in output.txt file

**Q5. Sol** – Following is the Comparison chart between computation rates of matrix matrix multiplication with and without blocking on **Local Machine**.

Matrix Multiplication with length 5000 took a lot of time to execute, hence it has not been taken. Instead I have implemented for length 1500, 2000, 2500



#### Output of the code:

Length :100 Matrix Mul Computation Rate: 3.13234e+08 block size:25 Block Matrix Mul Computation Rate: 335795836.1

Length :500 Matrix Mul Computation Rate: 489049210.1 block size:25 Block Matrix Mul Computation Rate: 563661000.7

Length :1000 Matrix Mul Computation Rate: 482134622.1 block size:25 Block Matrix Mul Computation Rate: 501433221.5

Length :1500 Matrix Mul Computation Rate: 471428005.7 block size:25 Block Matrix Mul Computation Rate: 418927530.6

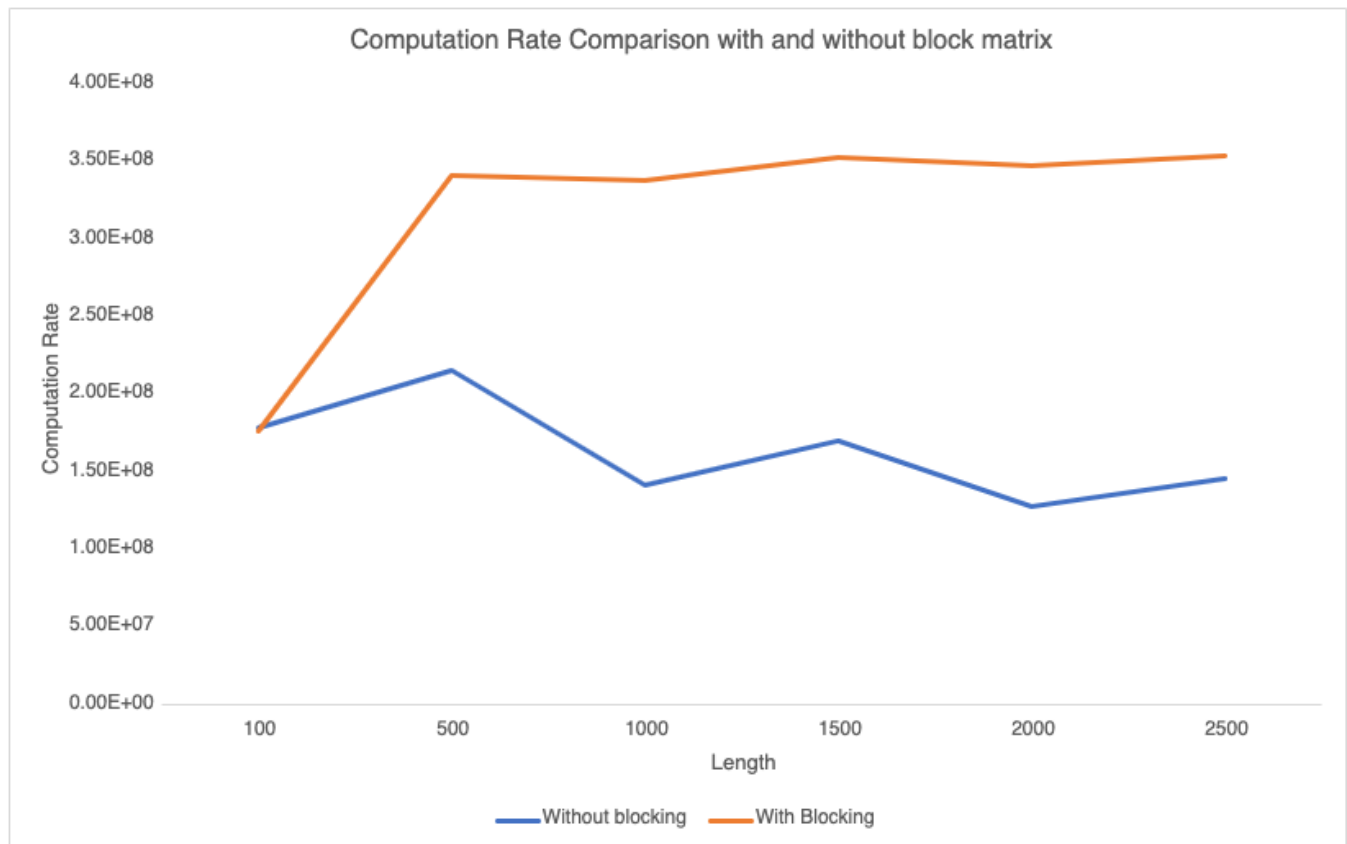
Length :2000 Matrix Mul Computation Rate: 188793389.6 block size:25 Block Matrix Mul Computation Rate: 424790283

Length :2500 Matrix Mul Computation Rate: 417884353.8 block size:25 Block Matrix Mul Computation Rate: 508297348.3



Following is the Comparison chart between computation rates of matrix matrix multiplication with and without blocking on **mc18 machine**

Matrix Multiplication with length 5000 took a lot of time to execute, hence it has not been taken. Instead I have implemented for length 1500, 2000, 2500



#### Output of the code:

Length :100 Matrix Mul Computation Rate: 1.78571e+08 block size:25 Block Matrix Mul Computation Rate: 176118351.5

Length :500 Matrix Mul Computation Rate: 215414548.1 block size:25 Block Matrix Mul Computation Rate: 340605951.6

Length :1000 Matrix Mul Computation Rate: 142068429.7 block size:25 Block Matrix Mul Computation Rate: 338242815

Length :1500 Matrix Mul Computation Rate: 170440101.1 block size:25 Block Matrix Mul Computation Rate: 352469889.7

Length :2000 Matrix Mul Computation Rate: 127862425.1 block size:25 Block Matrix Mul Computation Rate: 347065026.2

Length :2500 Matrix Mul Computation Rate: 146227473.9 block size:25 Block Matrix Mul Computation Rate: 353768984.6

**Execute the code:**

Code will be found inside hw1/question5/ folder

**Compile:** g++ main.cpp -o ./run

**Execute:** ./run

**Output:** output will be saved in output.txt file

From the outputs, we can comment that computation rate for block matrix multiplication is better than normal matrix multiplication as we can store the blocks of matrix in the cache based on the cache size and block size (here, the block size used is 25), leading to data reuse causing less cache misses and better performance as compared to iterating through entire matrix in normal matrix multiplication.