

SORT ON SINGLE SHARED MEMORY NODE

GOAL:

To gain experience programming with external data sort and multi-threaded programming.

DESIGN:

First we break the input data into different chunks and then sort each chunk using merge sort which has a time complexity of $O(n \log n)$, we then write them to different files named chunk1, chunk2, etc., and then perform a k-way merge on all the sorted chunks and then combine all of them together to get the sorted output file.

External Merge Sort and Linux Sort has been performed on the following files:

1. 1 GB
2. 4 GB
3. 16 GB
4. 64 GB

The input files have been generated by using gensort and the validation of the sorted files is performed using valsort. The memory usage of the shared memory sort has been limited to 8GB for the 16GB and 64GB datasets, while the smaller 1GB and 4 GB datasets have been sorted in the memory itself and then written to disk.

The Sorting has been implemented on a laptop with SSD to get the maximum speed and utilization. Multithreading while splitting the input file into different chunks has been implemented by extending the Thread class, where the classes which were created extends the java.Lang.Thread class.

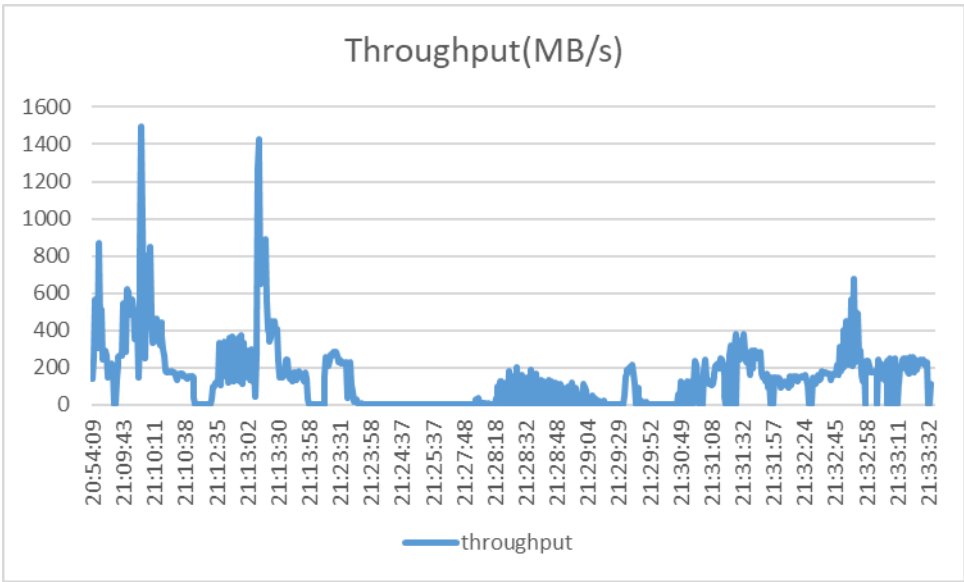
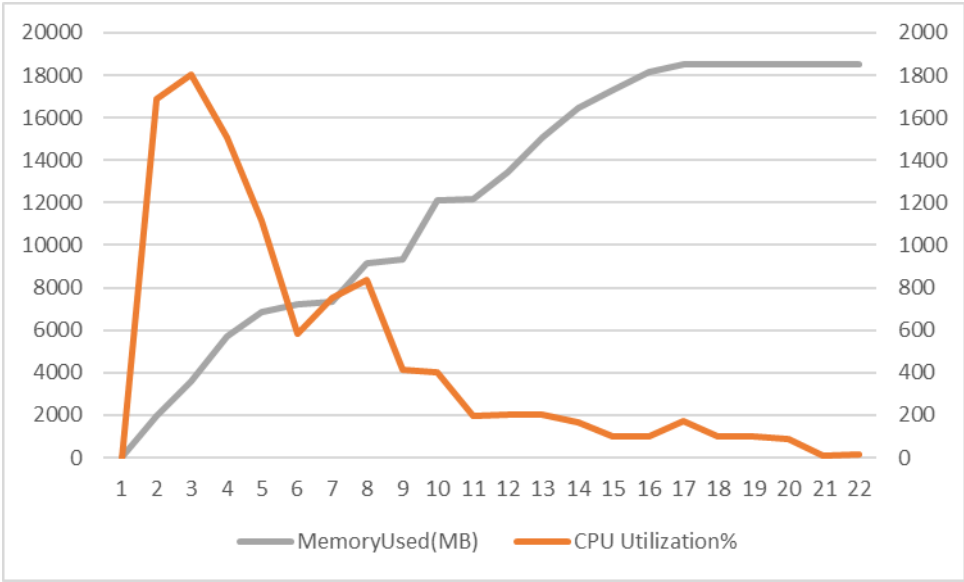
To make sure that the data is not cached in the OS memory before we run our experiments we execute the command “sudo sysctl -w vm.drop_caches=3” after every sort.

RESULTS:

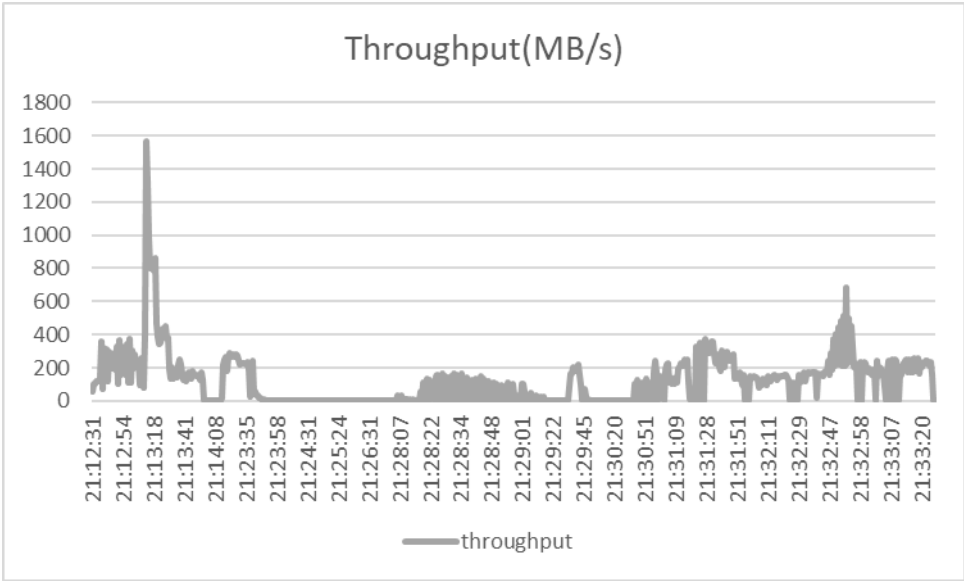
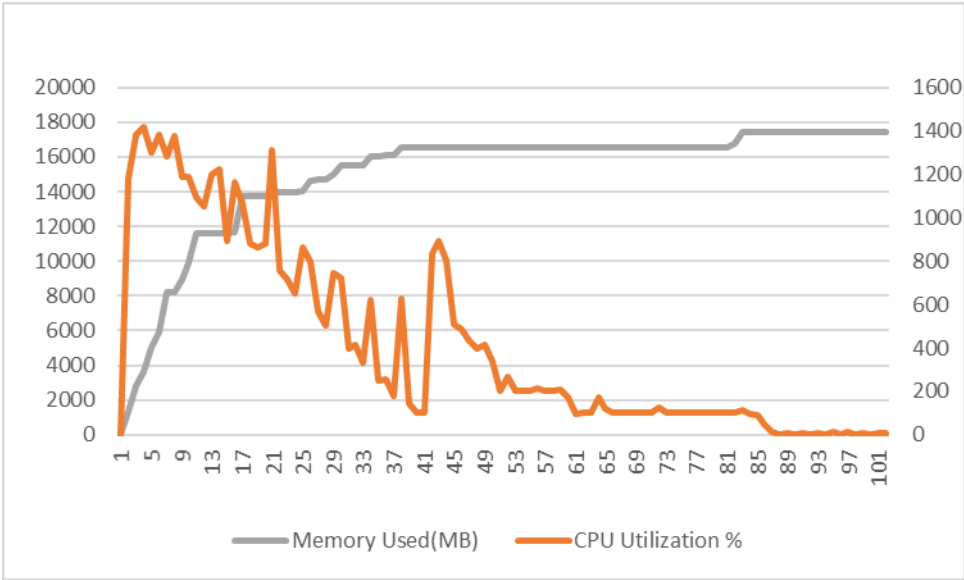
Experiment	Shared memory (1GB)	Linux sort (1GB)	Shared memory (4GB)	Linux sort (4GB)	Shared memory (16GB)	Linux sort (16GB)	Shared memory (64GB)	Linux sort (64GB)
Number of threads	1	1	8	8	16	16	16	16
Sort approach	Inside Memory	Inside Memory	External	External	External	External	External	External
Algorithm	Merge Sort	Merge Sort	Merge Sort	Merge Sort	Merge Sort	Merge Sort	Merge Sort	Merge Sort
Sort Time(seconds)	26	10	51	49	178	198	845	812
Data read(GB)	1	1	8	8	32	32	64	64
Data write(GB)	1	1	8	8	32	32	64	64
Overall, I/O throughput(MB/s)	145.94	180	119.93	228	68.03	261	179.8	297
Overall CPU Utilization(%)	4%	5%	14%	7%	16%	8%	17%	9%
Average memory Utilization(GB)	2457	1078	6954	4567	7856	6257	7658	7452

As we can see that Linux Sort outperforms MySort for 1GB,4GB and 64GB files, also it appears that Linux Sort appears to utilize the CPU and memory more efficiently than MySort. The Cloud Instance on which the Sorting was done has 48 cores and 96 threads while ps gives CPU utilisation as 100% for 1 full core utilisation, 200% for 2 full cores utilisation etc., the values were divided by 4800 to give utilisation as of the total cores that are available.

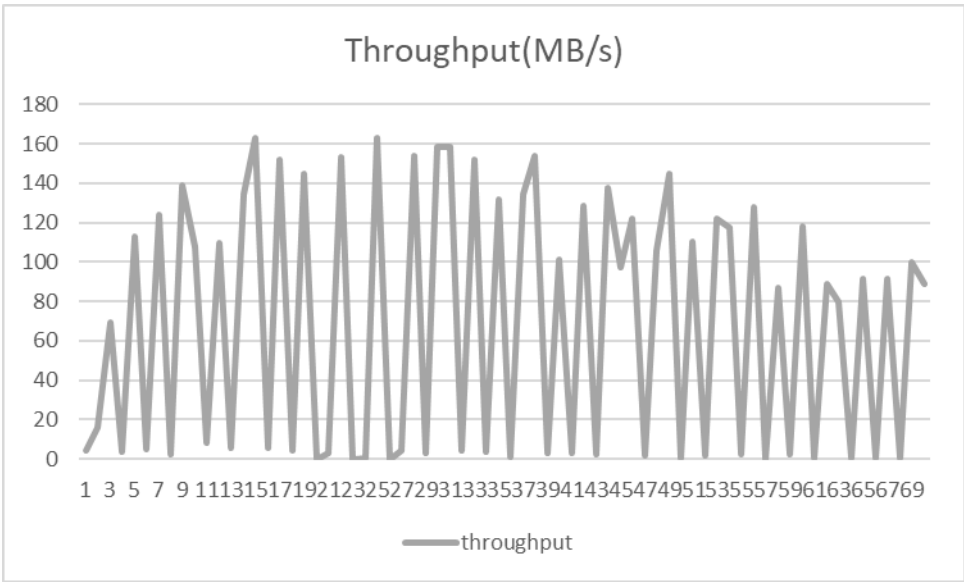
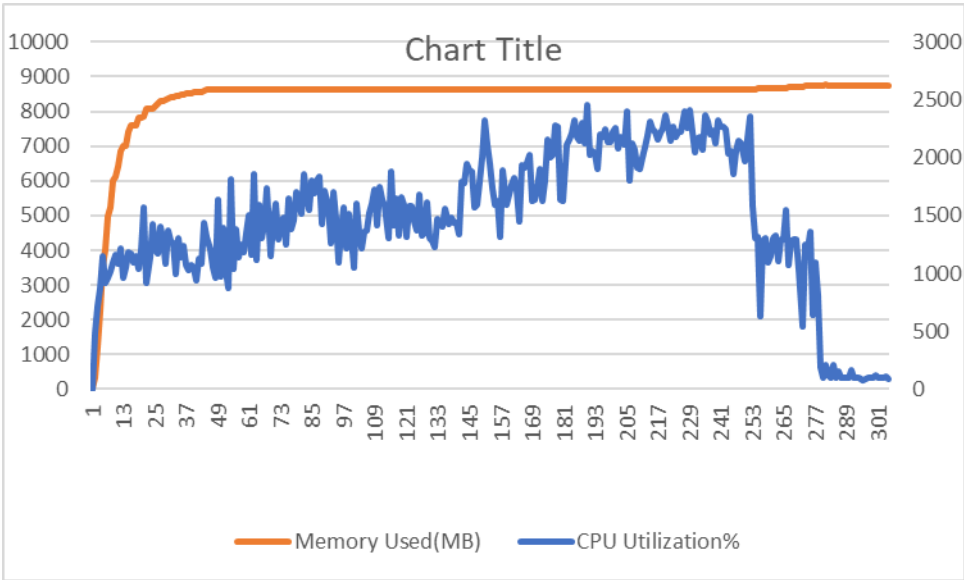
My Sort (1 GB):



My Sort (4 GB):



My Sort (16 GB):



My Sort (64 GB):

