# CS633 Assignment 1

*Naman Singla(200619)*
*Nishant Roshan(200643)*

## Different files and their functions

C Code(MPI Library) used to send data and measure time:

```c
Coding > CS633A > assign1 > C code.c > ⊗ main(int, char * [])
7    int main( int argc, char *argv[])
8    {
9        int myrank;
10       MPI_Status status;
11       double sTime, eTime, time, d;
12
13       MPI_Init(&argc, &argv);
14
15       int count = atoi (argv[1]);
16       char buf[count];
17
18       MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
19
20       // initialize data
21       for (int i=0; i<count; i++)
22           buf[i] = myrank+i;
23       MPI_Barrier(MPI_COMM_WORLD);
24       sTime = MPI_Wtime();
25       if (myrank == 0)
26           MPI_Send (buf, count, MPI_BYTE, 1, 1, MPI_COMM_WORLD);
27       else
28       if (myrank == 1)
29           MPI_Recv (buf, count, MPI_BYTE, 0, 1, MPI_COMM_WORLD, &status);
30       eTime = MPI_Wtime();
31       d = eTime-sTime;
32       MPI_Reduce (&d, &time, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);
33       if(myrank ==0 ){
34           printf ("%lf",time);
35       }
36       MPI_Finalize();
37       return 0;
38
39   }
```

Note:
- To measure the communication time, we make sure both processes have same starting time and then one ending last will decide the total communication time. This will calculate the precise total communication time that was required for transfer.
- Rank =0 is taken as root
- Run "ulimit -s unlimited" for sending large data.

Run scripts:

```bash
1    #!/bin/bash
2    echo "source,destination,size,time">"data.csv"
3    mpicc code.c
4    list=( $( ./available_nodes.sh ) )
5    len=${#list[@]}
6    for k in {1..3}
7    do
8        i=$(($RANDOM % $len))
9        j=$(($RANDOM % $len))
10       if [[ $i == $j ]]
11       then
12           j=$((($i+1)%$len))
13       fi
14       while read -r size
15       do
16               for z in {1..10}
17               do
18                   pairs="${list[$i]},${list[$j]}"
19                   t=$(mpirun -n 2 --hosts $pairs ./a.out $size </dev/null)
20                   echo "$pairs,$size,$t">>"data.csv"
21               done
22       done < "size.txt"
23   done
24
25   python3 plot.py
26
```

```bash
1    #!/bin/bash
2    for i in {1..50}
3    do
4        timeout 1 ping -c 1 -W 0.001 -s 10000 csews$i >/dev/null 2>&1
5        if [ $? == 0 ]
6        then
7            timeout 1 mpirun --hosts csews$i hostname >/dev/null 2>&1
8            if [ $? == 0 ]
9            then
10               echo "csews$i"
11           fi
12       fi
13   done
```

Note:
- available_nodes.sh produces a list that contains online csews pcs that can run mpi
- size.txt contains different sizes for which we want to acquire data.
- run.sh has following operations:
    - Compile the C code(code.c) with mpi.
    - Gets the list of online csews pcs
    - Selects 3 random pair of pc then execute code with different sizes for each pair 10 times.
    - File "data.csv" is generated which contains timings for different configurations.
    - Python script "plot.py" plots the data obtained and produces "plot.png"

```
Coding > CS633A > assign1 > 🐍 plot.py
  1    import pandas as pd
  2    import numpy as np
  3    import matplotlib.pyplot as plt
  4
  5    df = pd.read_csv("data.csv")
  6
  7    df.drop(["source","destination"],axis=1,inplace=True)
  8    df = df.astype({'size':'int','time':'float64'})
  9    # df['time']/=1000000
 10    plt.scatter(df['size'],df['time'],alpha=0.5)
 11    df = df.groupby('size')['time'].apply(np.hstack).to_frame().reset_index()
 12
 13
 14    xs=df['size']
 15    ys=df['time'].apply(lambda x: x.mean())
 16    yerr=[ys-df['time'].apply(lambda x: x.min()),df['time'].apply(lambda x: x.max())-ys]
 17    # print(yerr)
 18
 19    plt.plot(xs,ys)
 20    plt.errorbar(xs,ys,yerr=yerr,fmt ='o')
 21    plt.xlabel("Size(Bytes)")
 22    plt.ylabel("Time(s)")
 23    plt.savefig("plot.png")
```
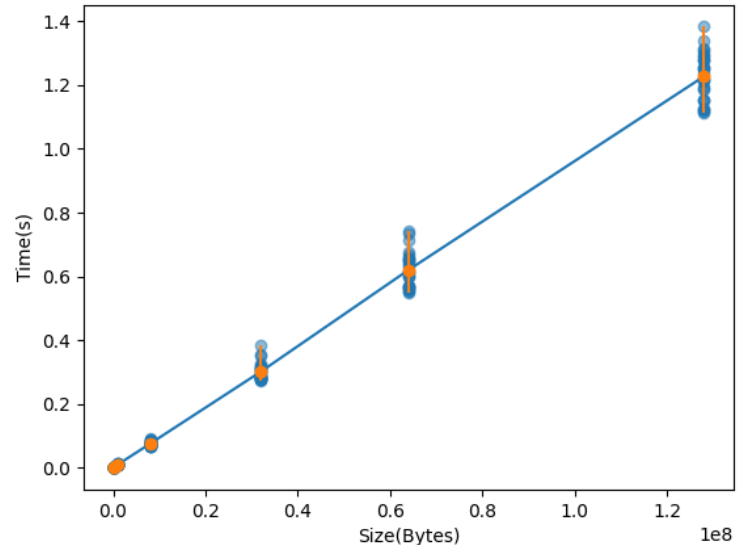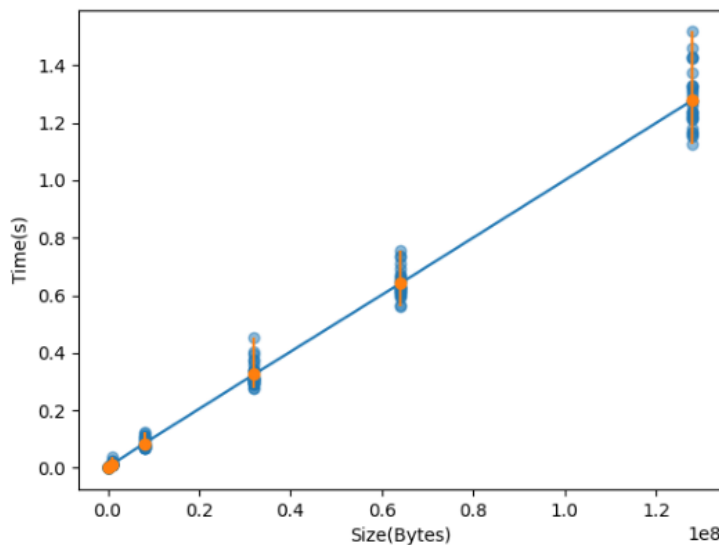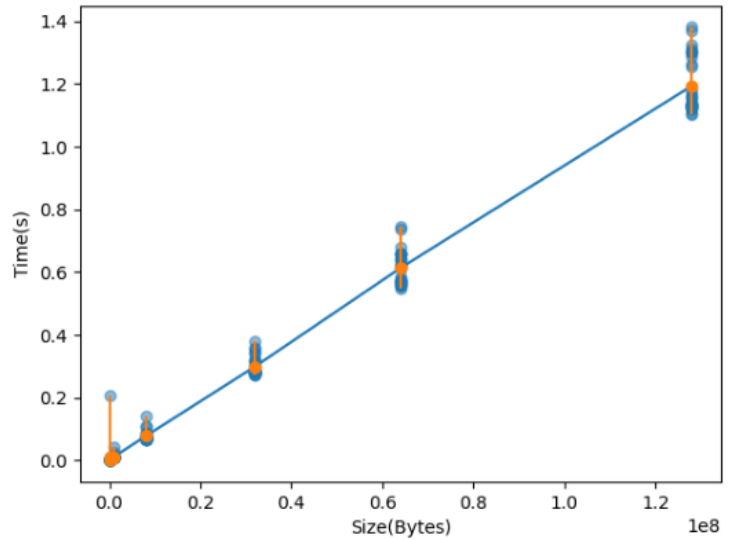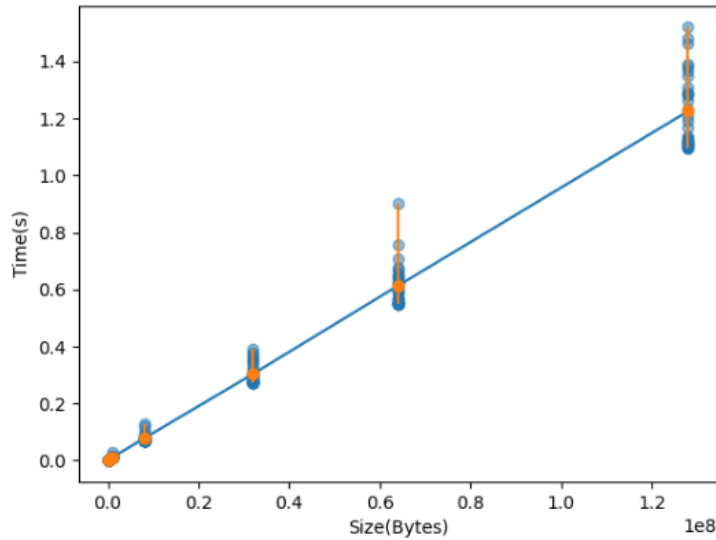
## How to run the submission?

- Make sure all the above-mentioned files are present in the current working directory.
- Run the bash script "run.sh"
- Output plot will be stored in files  plot.png

## Experimentation Methodology

The complete execution of the process is present in "run.sh".

1. Get the list of online nodes the have a reliable connection and can run mpi, using "available_nodes.sh"
2. Select a pair of nodes randomly from above obtained list.
3. Perform MPI_Send and MPI_Recv operations between above selected nodes.
   a. Put a barrier just before send and recv to make sure start time is same for both processes.
   b. Perform the transfer and note down the end time.
   c. Take the difference of  start time and the end time of process.
   d. Take maximum of above differences obtained for both processes.
   e. Value obtained above denotes the total time taken for communication, i.e both send and recv start at same time and the one ending last defines the total time taken. This time denotes that the complete data has been sent and received.
4. Perform Step 2,3 for each size mentioned in "size.txt". Take 10 readings for each size.
5. Perform 2-4 for 3 different pairs of nodes.
6. The data obtained in stored in csv file "data.csv".
7. We can plot the data obtained using matplotlib, numpy and pandas. The obtained plot is stored in "plot.png"

# Different plots obtained during experimentation and observation



1. The average time taken linearly varies with the size of the data sent.
   **Explanation:** This is evident if we assume a constant data rate. This may not be true for the small amount of data. But for large data sizes, data rates remain constant over the whole transfer span.
2. Error/deflections from expected behavior linearly increase with the amount of data sent(size).
   **Explanation:** Since data recorded for different sizes come from the same pair of nodes, we can say that different nodes may have different but constant data rates. Hence, we can claim that the deviation observed for a particular size can be scaled to get the approximate deviations of some other size value.
3. Upper deflection(+ve error) increases faster than lower deflection, i.e., the difference between the slope of max-time and mid-time is much more than that between the slopes of mid-time and min-time.
   **Explanation:** Maximum processes take the expected or near the expected completion time. Some processes experience high latency or lower bandwidth issues, so their time taken spikes high. But due to the number dominance average is only shifted slightly above the min.
4. In some cases, some outliers are observed(Their reading deflects by a large amount from the expected behavior).
   **Explanation:** These are cases when the daemon might be busy handling other users' requests, or sometimes the initialization cost goes high. These cases usually show a large increase in time. There is very less chance of an outlier being in lower deflection.
5. For each run slope of the line comes out to be approximately the same; only the error bars change.
   **Explanation:** The slope of the line denotes the data rate of transfer. Now, since all pcs are present on the same network using the same connections to the main network, we can say that their data rate is the same.