# Project 5 CS4473 report

To Ly
The University of Oklahoma
Norman, OK
tohly@ou.edu

**Problem 1:**

**Make plot, fitting and estimated values a linear regression function:**
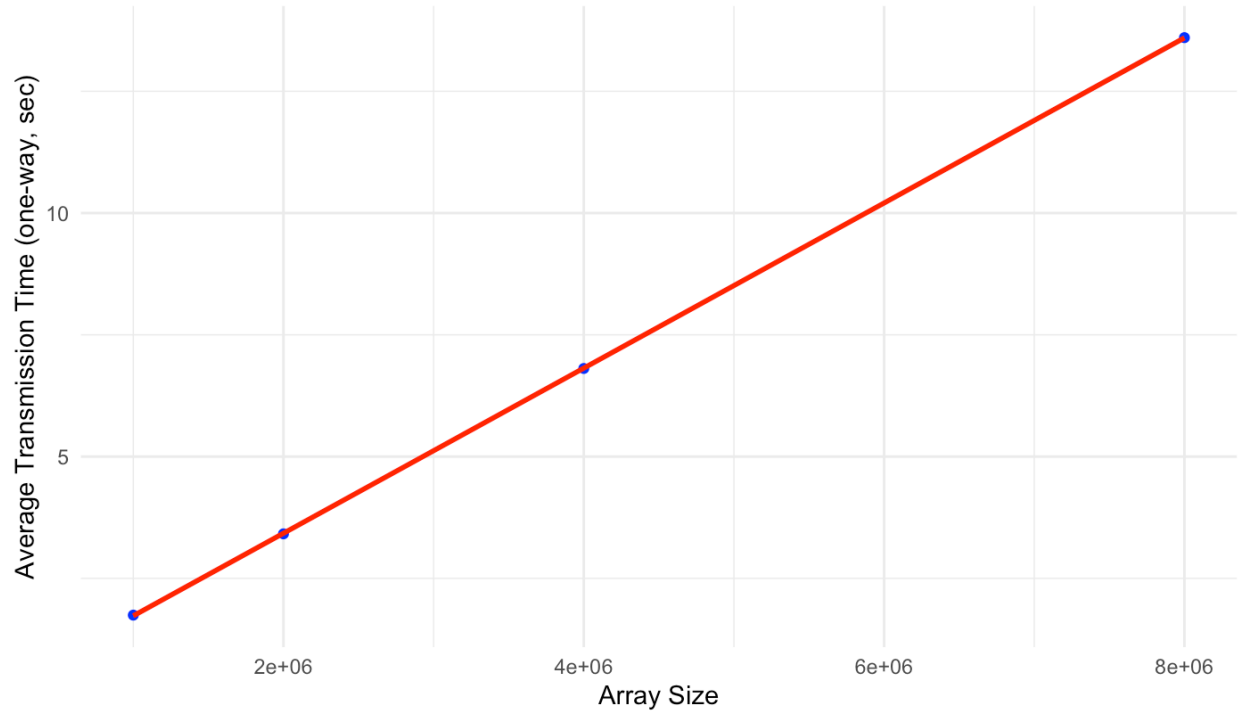
R-studio code:

```r
39   # Load necessary library
40   library(ggplot2)
41
42   # Read the data
43   diff_time_data <- read.csv("diff_time.csv")
44   same_time_data <- read.csv("same_time.csv")
45
46   # Define a function to perform linear regression and plot
47 ▾ plot_and_fit <- function(data, title) {
48     # Fit a linear model
49     model <- lm(time ~ array, data = data)
50
51     # Plot the data and the fitted line
52     plot <- ggplot(data, aes(x = array, y = time)) +
53       geom_point(color = "blue") +
54       geom_smooth(method = "lm", formula = y ~ x, color = "red") +
55       ggtitle(title) +
56       xlab("Array Size") +
57       ylab("Average Transmission Time (one-way, µs)") +
58       theme_minimal()
59
60     print(plot)
61
62     # Extract coefficients
63     coefficients <- coef(model)
64     latency <- coefficients[1] # Intercept
65     bandwidth_inverse <- coefficients[2] # Slope
66     bandwidth <- 1 / bandwidth_inverse
67
68     return(c("Latency (µs)" = latency, "Bandwidth (bytes/µs)" = bandwidth))
69 ▴ }
70
71   # Perform analysis for different nodes
72   results_diff <- plot_and_fit(diff_time_data, "Transmission Time vs. Array Size (Different Nodes)")
73
74   # Perform analysis for the same node
75   results_same <- plot_and_fit(same_time_data, "Transmission Time vs. Array Size (Same Node)")
76
77   # Print results
78   print(results_diff)
79   print(results_same)
80
```
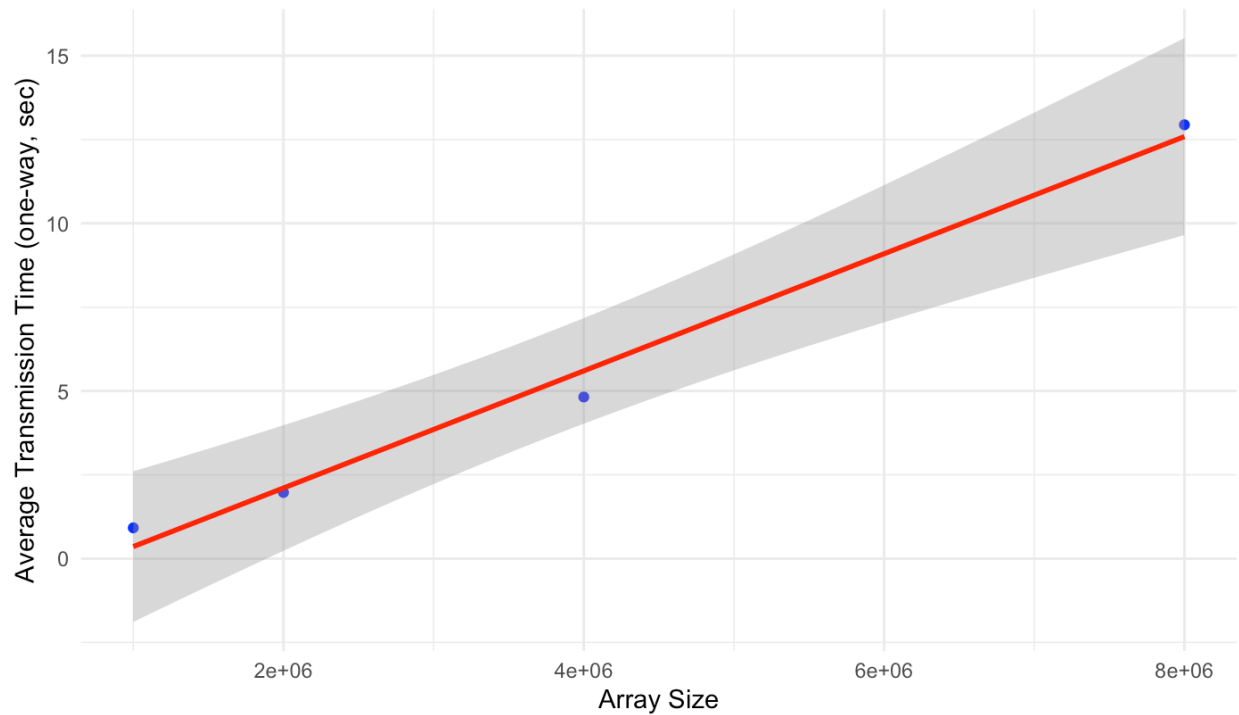
**Output:**
**Diff_node:**

Transmission Time vs. Array Size (Different Nodes)



**Same_node:**

Transmission Time vs. Array Size (Same Node)

**Estimate value:**

```
# Print results
print(results_diff)
```

```
##     Latency (sec).(Intercept) Bandwidth (bytes/sec).array
##                 3.626496e-02                 5.899792e+05
```

```
print(results_same)
```

```
##     Latency (sec).(Intercept) Bandwidth (bytes/sec).array
##                    -1.389076                 572422.789461
```

For the different nodes (two processes in two different nodes):
- Latency: Approximately 0.03626496 seconds
- Bandwidth: Approximately 589,979 bytes per second

For the same node (two processes on the same node):
- Latency: Approximately -1.389076 seconds (which is not physically meaningful and should be considered as an artifact of the model fitting, effectively it can be treated as zero)
- Bandwidth: Approximately 572,422 bytes per second

In this case, the same node run faster than diff node.

**Problem 2:**

| Array size | 262144 | 524288 | 1048576 |
|---|---|---|---|
| 2 processes | 0.000043775s | 0.000046349s | 0.000046137s |
| 4 processes | 0.000048155s | 0.000049294s | 0.000043162s |
| 8 processes | 0.000043579s | 0.000032312s | 0.000044462s |

In this case I run this program in same node sbatch.
We can see the processes increasing but the time taken to complete tasks does not decrease as expected for a scalable MPI program. Ideally, more processes should mean faster execution, but communication overhead, load imbalance, and the inherent non-parallelizable parts of the program seem to limit this speedup. In short, the MPI program exhibits limited scalability based on the times recorded for the different numbers of processes.

**Problem 3:**

| Array size | 262144 | 524288 | 1048576 |
|---|---|---|---|
| 4 processes on the same node | 0.020212s | 0.041939s | 0.088085s |
| 4 processes on 4 different nodes | 0.039590s | 0.076377s | 0.099020s |

Scalability on the Same Node:
The execution time increases with the array size, which is expected. However, the rate of increase seems sublinear, suggesting that the program scales quite well when adding more data within the same node. This could be due to the efficient use of shared resources like memory and cache, and avoiding the network overhead associated with communication between different nodes.

Scalability on Different Nodes:
Execution times are higher across the board compared to the same node scenario, which indicates network communication overhead. This overhead is particularly impactful in distributed systems where data must be sent over the network. The times increase with the array size, and the rate of increase seems slightly higher than within the same node, suggesting that network latency and bandwidth are becoming significant factors in the program's overall scalability.

The observation that 2 processes are slower than 4 suggests that the MPI program might have a high setup overhead or benefits from optimized communication with more processes. This implies that the program may not simply be limited by raw communication costs but also by how effectively it utilizes available resources and parallelizes tasks. Therefore, while 2-process communication seems less efficient, adding more processes improves performance, indicating a potential for better scalability when the program can leverage more parallelism and perhaps shared resources within a node.