# Project – COSC 6361.001

School of Engineering & Computer Sciences
*Texas A & M University at Corpus Christi*
Fall 2024

Points: 20

***Objective:***
Design and implement a parallel k-NN algorithm using the Message Passing Interface (MPI) to classify large datasets distributed across multiple processes. The core tasks will involve splitting the training dataset across multiple processors, performing distance calculations in parallel, and coordinating the selection of the nearest neighbors across all nodes.

**Team Members:**

1. Reddy Bhuvan Korlakunta (Team Leader)
2. Om Preetham Bandi
3. Thriveen Ullendula
4. Nandini Kodali
5. Sravyasri Virigineni
6. Sravani Jampani

**Table of Contents:**

---

## A. Introduction

The **k-Nearest Neighbors (k-NN)** algorithm is a simple yet effective supervised learning method widely used for classification and regression. However, its computational intensity, especially for large datasets, poses scalability and performance challenges due to the need to calculate distances for all training instances.

This project addresses these limitations by implementing a **parallel k-NN algorithm** using the **Message Passing Interface (MPI)** framework. By distributing the training dataset across multiple processes and parallelizing distance computations, the algorithm achieves improved efficiency, scalability, and reduced execution time.

The project involves dividing data, performing distributed calculations, and aggregating global results for classification. Experiments evaluated the algorithm's performance with varying kk-values, assessing accuracy, execution time, and scalability. The report highlights the design, implementation, and real-world applicability of the parallel k-NN, discussing its strengths, limitations, and potential for future optimization in parallel computing environments.

---

## B. Project Description

This project focuses on the implementation of a **parallel k-Nearest Neighbors (k-NN) algorithm** using the **Message Passing Interface (MPI)** framework. The k-NN algorithm is computationally intensive, as it requires calculating the distance between every test instance and all training instances. By leveraging parallel computing, this project aims to optimize the algorithm's execution time and scalability while maintaining its accuracy

## ❖ Algorithms:

**The Algorithms Used in Parallel k-NN**

The parallel implementation of the k-Nearest Neighbors (k-NN) algorithm is designed to optimize computation for large datasets by dividing the workload across multiple processes using the Message Passing Interface (MPI). Below are the key algorithms used, along with their explanations:

| Algorithm | Purpose | Key Operations |
|---|---|---|
| Data Distribution | Distribute training data among processes | Chunk allocation and indexing |
| Local Distance Computation | Compute distances locally | Vectorized distance calculations |
| Local k-NN Selection | Select k-nearest neighbors locally | Priority queue operations |
| Global k-NN Selection | Aggregate and select global k-nearest | Heap-based aggregation |
| Classification | Assign class labels to test instances | Weighted majority voting |
| Accuracy Evaluation | Measure classification accuracy | Label comparison and percentage calculation |

These algorithms collectively enable an efficient and scalable parallel k-NN implementation, demonstrating the power of MPI in handling large-scale data classification tasks.

---

## ❖ Local Data Structures for k-NN

To efficiently compute and manage the k-Nearest Neighbors (k-NN) in a parallel environment, the following local data structures are used:

I.  **Numpy Arrays for Training and Test Data:**

- **Code Location:** load_data() method
- **Purpose:** Training and test datasets are loaded as Numpy arrays for efficient computation.

```
load_data()

def load_data(self):
    if self.rank == 0:
        train_data = pd.read_csv(self.train_file, header=None).values
        test_data = pd.read_csv(self.test_file, header=None).values
        train_data = self.normalize_data(train_data)
        test_data = self.normalize_data(test_data)
    else:
        train_data = None
        test_data = None

    # Broadcast normalized datasets to all processes
    train_data = self.comm.bcast(train_data, root=0)
    test_data = self.comm.bcast(test_data, root=0)
    return train_data, test_data
```

II.     **Priority Queue (Heap) for Local k-NN Selection:**

- **Code Location:** find_local_knn() method
- **Purpose:** A priority queue (max-heap) is used to maintain the smallest k distances for a test instance.

```
find_local_knn()

def find_local_knn(self, local_train_data, test_instance, k):
    distances = self.compute_distances(local_train_data, test_instance)
    local_neighbors = [(distances[i], local_train_data[i, -1]) for i in
range(len(distances))]
    return heapq.nsmallest(k, local_neighbors, key=lambda x: x[0])
```

III.    **Distance Array for Local Distance Computation:**

- **Code Location:** compute_distances() method
- **Purpose:** Compute the squared Euclidean distance between the test instance and each training instance locally.

```
compute_distances()

def compute_distances(self, train_data, test_instance):
    squared_distances = np.sum((train_data[:, :-1] - test_instance[:-1]) ** 2,
axis=1)
    return squared_distances
```

IV.     **Local Neighbors List for Storing k-NNs:**

- **Code Location:** find_local_knn() method
- **Purpose:** Each process generates a list of tuples containing distances and labels of the local k-nearest neighbors.

```
compute_distances()

def compute_distances(self, train_data, test_instance):
    squared_distances = np.sum((train_data[:, :-1] - test_instance[:-1]) ** 2,
axis=1)
    return squared_distances
```

❖ **The Method of Selecting the Global k-NN**

### I.   Local k-NN Computation:

- **Code Location:** classify() method (calls find_local_knn() internally)
- **Purpose:** Each process finds the k-nearest neighbors for its local training dataset.

```
classify

for test_instance in test_data_iter:
    local_neighbors = self.find_local_knn(local_train_data, test_instance, k)
```

### II.   Gathering Local k-NNs:

- **Code Location:** gather_global_knn() method
- **Purpose:** Local neighbors from all processes are sent to the master process.

```
gather_global_knn

def gather_global_knn(self, local_neighbors, k):
    all_neighbors = self.comm.gather(local_neighbors, root=0)
```

### III.   Global k-NN Selection:

- **Code Location:** gather_global_knn() method
- **Purpose:** The master process combines the local results and selects the global k-nearest neighbors.

```
Global k-NN Selection

if self.rank == 0:
    combined_neighbors = [neighbor for sublist in all_neighbors for neighbor in
sublist]
    global_neighbors = heapq.nsmallest(k, combined_neighbors, key=lambda x: x[0])
    return global_neighbors
```

### IV.   Broadcasting Global k-NNs:

- **Code Location:** gather_global_knn() method
- **Purpose:** The global k-nearest neighbors are broadcasted back to all processes for final classification.

```
global_neighbors

global_neighbors = self.comm.bcast(global_neighbors, root=0)
```

**Key Advantages:**

- **Efficiency**: Local computation minimizes inter-process communication, reducing overhead.
- **Scalability**: The method scales well as it only exchanges the minimal necessary information (local k-nearest neighbors).
- **Accuracy**: Aggregating local results ensures the selection of the true global k-nearest neighbors.

These methods and data structures make the parallel k-NN algorithm both computationally efficient and scalable.

## ❖ The functionality of an algorithms

### 1. Data Distribution Algorithm

The training dataset is split among the available MPI processes to ensure an even workload. Each process is responsible for handling a portion of the training data. Steps**:**

1. Determine the chunk size for each process based on the total number of training samples and the number of processes.
2. Allocate extra samples to some processes if the dataset size is not perfectly divisible by the number of processes.
3. Share the test dataset with all processes so that distance calculations can be performed locally.

```
Data Distribution Algorithm

def distribute_training_data(self, train_data):
    num_samples = len(train_data)
    chunk_size = num_samples // self.size
    remainder = num_samples % self.size

    start_idx = self.rank * chunk_size + min(self.rank, remainder)
    end_idx = start_idx + chunk_size + (1 if self.rank < remainder else 0)
    local_train_data = train_data[start_idx:end_idx]
    return local_train_data
```

This algorithm ensures balanced data distribution among processes, minimizing idle time and maximizing parallel efficiency.

## 2. Local Distance Computation Algorithm

Each process computes the Euclidean distance between the test instance and the subset of training data it holds. This is the core computational step in k-NN.

1. For each test instance, subtract its feature vector from every training instance's feature vector in the process.
2. Square the differences, sum them across all dimensions, and compute the square root to get the Euclidean distance.

```
2. Local Distance Computation Algorithm

def compute_distances(self, train_data, test_instance):
    squared_distances = np.sum((train_data[:, :-1] - test_instance[:-1]) ** 2,
axis=1)
    return squared_distances
```

This algorithm uses a vectorized approach with NumPy for efficient computation, reducing the time complexity compared to a nested loop implementation.

## 3. Local k-Nearest Neighbors Selection Algorithm

Each process maintains a priority queue (max-heap) to identify the k-nearest neighbors within its local dataset.

1. Compute the distance between the test instance and all training instances in the process.
2. Use a heap to store the k smallest distances along with their corresponding labels.
3. Ensure efficient insertion and removal operations using the heap structure.

```
3. Local k-Nearest Neighbors Selection Algorithm

def find_local_knn(self, local_train_data, test_instance, k):
    distances = self.compute_distances(local_train_data, test_instance)
    local_neighbors = [(distances[i], local_train_data[i, -1]) for i in
range(len(distances))]
    return heapq.nsmallest(k, local_neighbors, key=lambda x: x[0])
```

This algorithm efficiently selects k-nearest neighbors from the local subset, reducing memory usage and computational overhead.

## 4. Global k-Nearest Neighbors Selection Algorithm

The global k-nearest neighbors are identified by aggregating the results from all processes.

1. Each process sends its local k-nearest neighbors to the master process.
2. The master process combines all the local results.
3. A max-heap is used to identify the global k-nearest neighbors from the combined results.

```python
4. Global k-Nearest Neighbors Selection Algorithm

def gather_global_knn(self, local_neighbors, k):
    all_neighbors = self.comm.gather(local_neighbors, root=0)
    if self.rank == 0:
        combined_neighbors = [neighbor for sublist in all_neighbors for neighbor in
sublist]
        global_neighbors = heapq.nsmallest(k, combined_neighbors, key=lambda x:
x[0])
        return global_neighbors
    return None
```

This step minimizes inter-process communication overhead by sending only the top-k results from each process. The use of a heap ensures efficient selection of the global neighbors.

## 5. Classification Algorithm

The class label for each test instance is determined using a majority vote among the global k-nearest neighbors.

1. Count the occurrences of each class label among the k-nearest neighbors.
2. Use a weighted voting scheme, where weights are the inverse of the distances, to resolve ties or improve accuracy.
3. Assign the class label with the highest weighted count.

```python
5. Classification Algorithm

def classify_instance(self, global_neighbors):
    weights = {}
    for distance, label in global_neighbors:
        weights[label] = weights.get(label, 0) + 1 / (distance + 1e-5)  # Weighted
by inverse distance
    return max(weights, key=weights.get)
```

This algorithm ensures robust classification by incorporating distance-based weights, making it less susceptible to noise or outliers.

**6. Accuracy Evaluation Algorithm**

The accuracy of the classifier is evaluated by comparing predicted labels to the true labels in the test dataset.

1. Compute the proportion of correctly classified test instances.
2. Express accuracy as a percentage.

```
6. Accuracy Evaluation Algorithm

def evaluate(self, predictions, test_data):
    true_labels = test_data[:, -1].astype(int)
    accuracy = np.mean(np.array(predictions) == true_labels) * 100
    return accuracy
```

This algorithm provides a quantitative measure of the classifier's performance, enabling comparison across different configurations.

❖ **The Running Time Complexity of an algorithm**

The running time complexity of the parallel k-NN algorithm is divided into three main phases: **data distribution**, **local distance computation**, and **global aggregation**.

1. **Data Distribution:**

   The training data is split across PP processes, and the test dataset is broadcasted to all processes.
   ➢ This step has a complexity of $O(ntest \times m)$ for broadcasting.
   ➢ The complexity for distributing the training data is $O(ntrain / P)$), where ntrain is the number of training samples, ntest is the number of test samples, and mm is the feature dimension.

2. **Local Distance Computation:**

   Each process computes distances between its local training data and the test instances.
   ➢ This step dominates the computation, with a complexity of $O((ntrain / P) \times ntest \times m)$.

3. **Global Aggregation:**

   Each process sends its k-nearest neighbors to the master process.
   ➢ The communication cost is $O(P \times k)$.
   ➢ The master process aggregates results and selects global k-nearest neighbors in $O(P \times k \times logk)$.

Overall, the parallelization reduces computational load per process while adding minimal communication overhead. Compared to sequential k-NN, the parallel approach improves efficiency and scalability for large datasets by distributing the workload among PP processes.

## ❖ The working complexity of an algorithm

The **working complexity** of an algorithm measures the total computational effort required to complete a task across all processes in a parallel system. It accounts for both local computations performed by each process and the global communication and aggregation overhead.

In the parallel k-Nearest Neighbors (k-NN) algorithm, the working complexity comprises three main steps:

    **1.**     **Local Distance Computation:** Each process calculates the distances between its portion of the training data and the test instances. Across all processes, the total complexity is $O(n \cdot m)$, where n is the total number of training samples, and mm is the number of features.

    **2.**     **Local k-NN Selection:** Each process uses a heap to find the kk-nearest neighbors from its local training subset. The total complexity across all processes is $O(n \cdot \log k)$, where kk is the number of neighbors.

    **3.**     **Global k-NN Selection:** The master process aggregates the local neighbors and selects the global k-nearest neighbors using a heap. The complexity of this step is $O(P \cdot k \cdot \log k)$, where PP is the number of processes.

The overall working complexity of the algorithm is: $\mathbf{O(T \cdot n \cdot m + T \cdot n \cdot \log k + T \cdot P \cdot k \cdot \log k)}$

where T is the number of test instances. This ensures scalability and efficiency in distributed environments.

---

## ❖ The Speedup Complexity of an Algorithm

The **speedup complexity** of an algorithm shows how much faster it runs in parallel compared to sequential execution, **Sequential Execution Complexity**: In sequential k-NN,

The time complexity is: $\mathbf{O(T \cdot n \cdot m + T \cdot n \cdot \log k)}$

where TT is the number of test instances, n is the number of training samples, mm is the feature size, and kk is the number of neighbors.

1. **Parallel Execution Complexity**:
   - When n training samples are divided across PP processes,
   - The time per process is: $\mathbf{O(T \cdot \frac{n}{P} \cdot m + T \cdot \frac{n}{P} \cdot \log k)}$

2. **Speedup**:
   - Without communication overhead: Speedup is approximately $O(P)$ (ideal linear speedup).
   - With communication: Speedup is slightly less than $O(P)$, depending on the dataset size and the number of processes. For large n, speedup remains close to linear.

---

### ❖ The Efficiency of an Algorithm

The **efficiency** of an algorithm is a measure of how well it utilizes computational resources in a parallel environment. It is calculated as the ratio of speedup to the number of processes used. If all processes are effectively contributing to reducing execution time, the efficiency is close to 1 (or 100%). In the parallel k-NN algorithm:

- **High Efficiency**: Achieved when the dataset is large, and the workload is evenly distributed across processes.
- **Low Efficiency**: Balancing data distribution and reducing communication between processes are critical to maintaining high efficiency.

### ❖ Scalability of an Algorithm

**Scalability** measures how well an algorithm performs as the problem size or the number of processes increases. It is divided into two types:

1. **Strong Scalability**:
   - The algorithm is tested with a fixed problem size while increasing the number of processes.
   - In an ideal scenario, execution time decreases proportionally to the number of processes. However, as more processes are added, communication overhead might limit scalability.
2. **Weak Scalability**:
   - A scalable algorithm maintains consistent execution time as both the dataset and number of processes grow. The parallel k-NN algorithm demonstrates good weak scalability due to its ability to handle larger datasets efficiently.

### ❖ Experiments and Results Discussions

The results of running the **parallel k-Nearest Neighbors (k-NN)** algorithm using MPI. Below is a breakdown of the displayed terminal output:

## 1. Initialization

- **Number of Processes:**
  - Number of Processes: 4 indicates that the parallel computation is distributed across 4 MPI processes.
  - Each process is responsible for a subset of the training data.
- **Data Loading:**
  - Process 0 handles data loading and shows:
  - **Training Data Shape:** (8192, 16) indicates 8192 samples with 16 features each.
  - **Test Data Shape:** (1024, 16) indicates 1024 test samples, each with 16 features.

## 2. Training Data Distribution

- The training dataset is split evenly across 4 processes.
- **Training Data Chunk Sizes:**
  - Each process gets 2048 training samples with 16 features (total training samples ÷ number of processes => 8192÷4=2048).

## 3. Processing Test Instances

- The system processes the test dataset of 1024 samples.
- **Processing 1024 test instances** confirms all test samples are being evaluated.
- **Processed 1024/1024 test instances** indicates that all test samples have been classified successfully.

## 4. Results for k-Nearest Neighbors

The table shows the results for different $kk$-values (number of nearest neighbors):

- **k = 8**:
  - **Accuracy:** 85.55%
  - **Execution Time:** 0.4536 seconds
  - Observation: A small $kk$-value results in relatively fast execution and good accuracy.
- **k = 16**:
  - **Accuracy:** 86.13%
  - **Execution Time:** 0.4745 seconds
  - Observation: Increasing $kk$ slightly improves accuracy with a minimal increase in execution time.
- **k = 32**:
  - **Accuracy:** 85.64%
  - **Execution Time:** 0.5615 seconds
  - Observation: Larger $kk$ slightly reduces accuracy while increasing computation time.

## 5. Summary

- The algorithm provides classification results for all test instances, varying $kk$-values to evaluate performance. Execution times increase as $kk$ grows due to more neighbors being evaluated and communicated.
- Accuracy remains relatively stable across different $kk$-values, with the highest accuracy observed at k=16.

**Final Observations**

- **Efficiency:** The data distribution and parallel computation ensure efficient processing with relatively low execution times.
- **Scalability:** The output demonstrates strong scalability as all processes contribute evenly.
- **Customization:** The input prompt for kk-values allows flexibility for experimenting with different settings.

This output validates the correctness and performance of the parallel k-NN implementation in a distributed computing environment. Let me know if you'd like more details or further analysis.

**Pros and Cons Discussions**

**Pros**:

1. **Improved Performance**: Parallelization significantly reduces execution time for large datasets.
2. **Scalability**: The algorithm handles increase in dataset size well, making it suitable for real-world applications.
3. **Flexibility**: The use of MPI allows the algorithm to be deployed on various distributed systems, from local clusters to high-performance computing environments.

**Cons**:

1. **Communication Overhead**: Global aggregation of k-nearest neighbors can become a bottleneck, especially for small datasets or large kk.
2. **Imbalanced Workload**: Uneven distribution of training data can lead to some processes being underutilized.
3. **Memory Usage**: Each process requires a copy of the test dataset, which may become problematic for very large datasets.

**Further Improvements**

1. **Load Balancing**:
   - Implement dynamic data partitioning to ensure all processes handle an equal workload.
   - Use workload profiling to predict and adjust imbalances in real time.
2. **Optimization of Communication**:
   - Reduce communication overhead by using more efficient MPI collective operations.
   - Compress or aggregate data before communication to minimize the size of transmitted messages.
3. **Alternative Distance Metrics**:
   - Support for non-Euclidean distance metrics (e.g., Manhattan, cosine similarity) could improve accuracy for specific datasets.
4. **GPU Acceleration**:
   - Leverage GPUs for local distance computations to further speed up the algorithm, especially for datasets with high-dimensional features.

These improvements would enhance the parallel k-NN algorithm's efficiency, scalability, and applicability to more complex datasets and environments.