# A Combined Approach For Clustering Biological Ligand Using Hadoop and MPI

[1]Swagita Dwivedi, [2]Lalitkumar P.Bhaiya

Dept. of Computer Science and Engineering,
Bharti College of Engineering, Chhattisgarh, India

*Abstract :* Driven by the steady extended in "conceived advanced" data, data driven registering gives new troubles in data organization and investigation in the life sciences and in various controls stretching out from coordination's (movement examination) to sociology(online organizing). These troubles are through and through not quite the same as the standard reenactment challenges since dealing with an extensive number of individual records and terabytes of data. Cutting edge investigation strategies in light of packing and requests for ordinary reenactments consistently require the examination of complex data records in an iterative methodology. In this paper we present MPI-MapReduce based clustering of complex biological dataset utilizing octree clustering technique. We performed investigates a conveyed environment utilizing MapReduce and MPI. The outcomes beats the current clustering technique which utilizes Euclidian distance. The proposed L2 standard based octree clustering demonstrates more levels of isolated ligand**.**

*IndexTerms* **- Biological Ligand, MPI, MapReduce, Octree Based Clustering.**

## I. INTRODUCTION

As of late, novel patterns towards parallel structures have offered better response time and execution for parallel calculation applications. Multi-center clusters show noticeable changes in equipment structures in the midst of HPC ages, for instance, symmetric multi-processors (SMP shared memory systems) or single center hubs. Multi-center gatherings impact this feasible to mix two features of parallel programming languages; to message going in shared memory and multi-cores using shared memory.

Executing parallel estimation brought various thoughts which need to think about just, memory progressive system [1]. It will impact execution by two critical parameters: memory inertness; that is the time passed from point a touch of data is required until the point that data wind up evidently available, memory transfer speed; that is the speed which data are sent from memory to processors.

This investigation focused on how multi-center packs designs can impact the execution of message passing and shared memory while running parallel calculation application, and the ways to deal with upgrade taking care of time and execution by mix of two essential parallel programming lingos: unadulterated MPI and OpenMP.

This investigation focused on how multi-center gathering's designs can impact the execution of message passing and shared memory while running parallel calculation application, and the ways to deal with improve taking care of time and execution by mix of two essential parallel programming language: unadulterated MPI and OpenMP.

## II. MESSAGE PASSING INTERFACE

MPI outfits parallel equipment dealers with an evidently described base course of action of routines that can be successfully executed [2]. Hence, equipment dealers can develop this gathering of standard low-level routines to make increasingly raised sum routines for the distributed memory correspondence condition furnished with their parallel machines. MPI gives a simple to-use advantageous interface for the basic customer, yet one adequately extreme to empower software engineers to use the prevalent message passing activities available on cutting edge machines.

### 2.1 Functionality

The MPI interface is proposed to give fundamental virtual topology, synchronization, and correspondence value between a course of action of techniques (that have been mapped to hubs/servers/PC events) in a tongue autonomous manner, with language specific semantic structure (ties), notwithstanding a couple of language specific features.

### 2.2 Concepts

MPI gives a rich extent of capacities. Some of them are introduced below.

1. Communicator: Communicator objects interface gatherings of methods in the MPI session. Each communicator gives each contained system a free identifier and organizes its contained methodology in a mentioned topology.

2. Point-to-point basics: Various basic MPI capacities incorporate correspondence between two specific systems. A predominant case is MPI Send, which empowers one demonstrated strategy to establish a connection on a moment determined procedure.

3. Aggregate fundamentals: Aggregate capacities incorporate correspondence among all systems in a strategy gathering (which can mean the entire methodology pool or a program-portrayed subset). An ordinary work is the MPI_Bcast call (another approach to state "communicate").

## III. HADOOP – MAPREDUCE

The Apache Hadoop is a system that considers the conveyed handling of broad educational accumulations across over gatherings of PCs using fundamental programming models [3]. It is proposed to scale up from single servers to countless, each offering neighborhood computation and capacity [4]. The essential mapreduce task are introduced in fig. 1.

Hadoop engages the application to work in a distributed environment [5]. There may be a large number conveyed part participating to complete a single errand. Generally, the huge log records are appropriated over various clusters known as HDFS cluster [12] (Hadoop distributed file system). Hadoop isolates the records into the amount of pieces. These pieces are appropriated over various clusters and dealt with in each structure in a parallel plan. The execution of the Hadoop structure is gotten by working archives in a parallel environment.
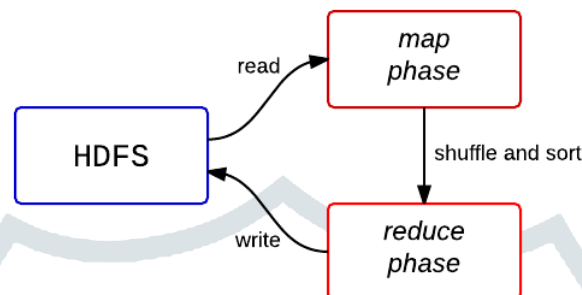


Fig. 1. Shows essential activity of MapReduce Model

## IV. RELATED WORK

K. A. Abdul Nazeer et al. [6] proposes k-means algorithm, for assortment of arrangements of values of starting centroids, produces various clusters. Result cluster quality in algorithm relies upon the choosing of initial centroids. Two stages consolidate into genuine k means algorithm: first to deciding introductory centroids and second to designating information points to the nearest clusters and after that re-figuring the clustering mean.

Soumi Ghosh et al. [7] present a general trade of two clustering algorithms in specific centroid based K-Means and representative object based FCM (Fuzzy C-Means) clustering algorithms. This discourse is on the basis of assessment of execution of the effectiveness of clustering yield by applying these algorithms.

Shafeeq et al. [8] present a changed K-means algorithm to improve the cluster quality and to fix the optimal numerous amounts of groups. As information number of clusters (K) given to the K-means algorithm by the client. In any case, in the useful situation, it is remarkably basic to fix the numerous amounts of clusters early. The method proposed in this paper works for both the cases for instance for known number of clusters early similarly as obscure number of clusters. The client has the flexibility either to fix the quantity of clusters or information the base number of clusters required. The new cluster centers are figured by the algorithm by growing the cluster counter by one in each cycle until it fulfills the authenticity of cluster quality. This algorithm will beat this issue by finding the satisfy less measure of clusters on the run.

Junatao Wang et al. [9], in this paper creator propose an improved k-means algorithm using noise information filter. The deficiencies of the ordinary k-means clustering algorithm are overwhelmed by this proposed algorithm. The algorithm makes density set up together detection strategies based with respect to features of noise information where the revelation and processing ventures of the noise information are appended to the first algorithm. By pre-processing the information to dismiss these noisy data before gathering informational collections the cluster union of the clustering yield is improved inside and out and the effect of noise information on k-means algorithm is diminished enough and the clustering yield are progressively exact.

Shi Na et al. [10] present the analysis of shortcomings of the improved k-means algorithm. As k-means algorithm needs to ascertain the distance between each datum object and all cluster focuses in each iteration. This recurrence process impacts the productivity of clustering algorithm. An upgraded k-means algorithm is proposed in this paper. A basic information structure is expected to store some information in every iteration which is to be used in the following iteration. Computation of distance in every iteration is stayed far from by the proposed system and saves the running time.

## V. PROBLEM IDENTIFICATION

In investigations of disease frames, a run of the mill issue is to search for little molecules (ligands) that can collaborate with a greater molecule, for instance, a protein when the protein is locked in with a disease state. All the more explicitly, when ligands dock well in a protein, they can possibly be used as a medicine to stop or deflect diseases related with the protein's glitch. This is a normal sort of investigation done in cure headway look at. In this paper, for example of its propriety, we apply our method to fundamental science datasets containing considerable amounts of ligand adjustments. This paper means to group all the related qualities through clustering and MapReduce algorithm which adroitly chooses the ligands of comparable kinds.

Usually, an innocent methodology used to group comparable basic science compliances is through geometry-based clustering. Such technique necessitates that data is secured in a moved zone remembering the ultimate objective to apply the Root-Mean Square Deviation of each ligand with the different ligands in the dataset. The investigation requires the nuclear dataset to be moved entire into a focal server. This totally thought methodology isn't versatile and can achieve some genuine stockpiling and transfer speed issue on the server side.

## VI. METHODOLOGY

In this area we present the proposed procedure for actualizing octree based clustering utilizing conveyed vector distance estimation with Map Reduce and MPI.

We proposed a novel instrument for investigating extensive measure of information in brief timeframe with highest accuracy. We have thought about biological qualities information. The analyses are performed on an distributed environment. The proposed strategy endeavors to beat the current system in which the information development occurred for preparing immense volume of information.

Proposed component insightfully chooses little arrangement of information and transfer among hubs and group those utilizing octree based clustering instrument. The real advance contains an information reshaping that is associated with each individual information record at the same time. This movement adventures relevant isolated information properties into a sort of metadata. We execute a one of kind techniques to supervise and look at the reshaped information.

In the main variety, called Metadata Movement, the removed properties are moved transversely over hubs to gather an overall viewpoint of the dataset; these properties are iteratively and secretly separated while chasing down some class or cluster association. The proposed technique relies upon the general subject of moving estimation to the information. We join our algorithm into the MapReduce perspective [2][11] and use them to consider a complex essential science dataset of a large number of ligand particles.

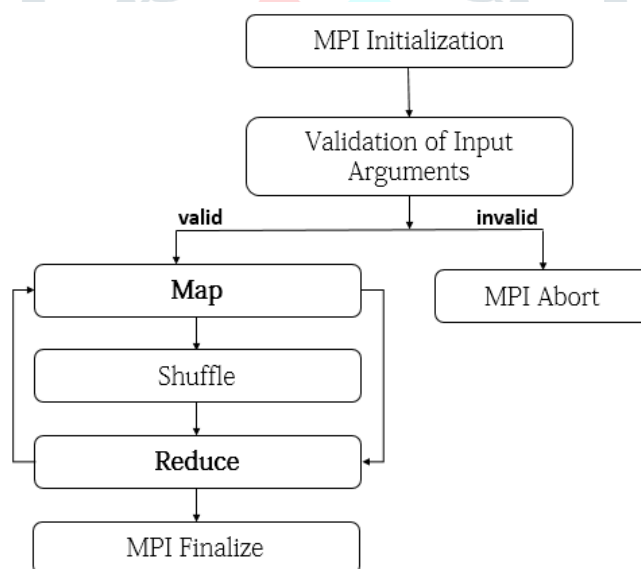The fig. 2. Demonstrates the proposed strategy step by step.



*Fig. 2.* Essential Workflow of MPI-MapReduce Clustering Framework

**Step 01: Initialization of MPI Parameters**

The MPI Parameters are set up for execution in parallel environment. The function, for example:

1. MPI_INIT()
2. MPI_Comm_Rank()
3. MPI_Comm_Size()

**Step 02: Map Phase**

Map stage utilized for octree key for every ligand.

Key → Octree Key
Value → Id for each ligand

It computes the octree key for every ligand by taking contribution from the biological information which is the gathering of different molecules or ligands.

**Step 03: Shuffling Phase**

Now, the key, esteem sets are divided among the r reduce forms dependent on the principal couple of digits of the key. Next the reducer perform local decrease, and chooses which level to go straightaway. The distance for every one of the properties are determined here and passed to the reducer for various level handling.

**Step 04: Reducing Phase**

The reducer principle task is to perform decrease of information and selecting best level for partitioning of information focuses. It iterated through entire archive and parallel group every one of them dependent on local properties. At long last, it consolidates result and put away into documents with level data.

*Properties Exchange*

Given a ligand with p nuclear directions (xi, yi, zi, with I from 1 to p), we play out a projection of the directions in their particular 3 planes (x, y), (y, z), and (z, x). Every projection results in a lot of 2D focuses on the related 2D plane.

Extracted properties can be capriciously disseminated over the 3-dimensional space of property focuses and over the hubs of the conveyed memory system.

In Metadata movement, the extracted properties are exchanged among hubs over the conveyed memory system to reconstruct a global perspective on the data content in the logical information with the goal that comparative properties in the end live on a similar hub or hubs that are topologically near one another as appeared in fig. 3.
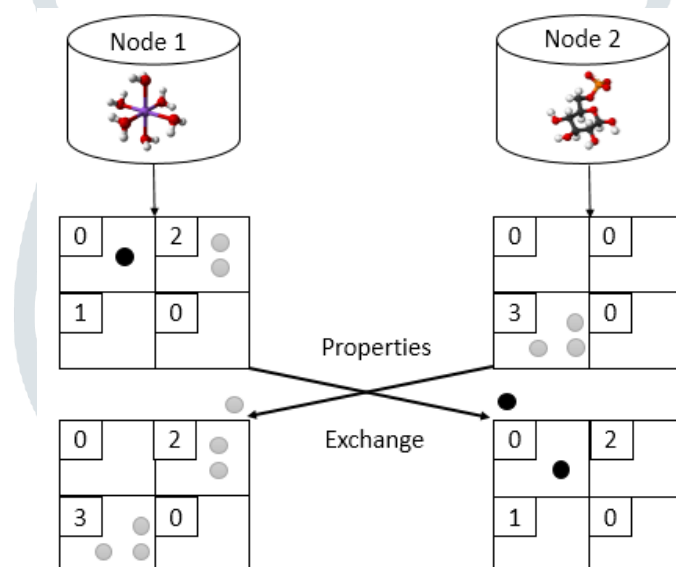


*Fig. 3.* Demonstrates the Exchange of metadata properties *in nodes*

**VII. RESULTS AND DISCUSSIONS**

We have performed experiments utilizing Biological dataset which comprises of 76,889 ligand molecules. The system configuration is introduced in Table I.

*Table I. System Configuration*

| S.NO | Attributes | Values |
|------|------------|--------|
| 1 | Software Used | Open-MPI MapReduce |
| 2 | Language Used | C++ |
| 3 | Algorithm | Octree Based Clustering, L2 Norm Distance Function |
| 4 | Dataset | Biological Dataset |

Utilizing following system configuration the clustering of ligands are finished. For clustering we have considered different threshold values, for example, 300, 500 and 1000 incentive for clustering of molecules. The threshold used to restrict the output of the group to contain explicit scopes of molecules and density.

We have performed experiment while thinking about Euclidian distance (existing strategy) and L2 Norm for disseminated calculation (proposed technique) for examination.

The execution time and clustered records are appeared Table II.

Table II. Clustered Records and Convergence for Execution Comparing Euclidian and L2 Norm Calculation

| Threshold | Convergence Levels | Clustered Records | Level Information |
|---|---|---|---|
| 300_eu | 4 | 71 | Level 4 Only |
| 300_l2 | 9 | 64 | Level 8 and 9 |
| 500_eu | 4 | 31 | Level 4 Only |
| 500_l2 | 9 | 30 | Level 8 and 9 |
| 1000_eu | 4 | 3 | Level 4 Only |
| 1000_l2 | 9 | 11 | Level 8 and 9 |

The results gives clear sign that, for threshold estimation of 300 with Euclidian distance the convergence rate is 4 and records clustered are 71 however with L2 norm its combined at level 9 with 64 information records. Be that as it may, the information are consistently circulated all through level 8 and 9. Not at all like with Euclidian distance are the information focuses available at level 4 as it were. In any case, in actuality the biological information indicates are be disseminated some place around level 8 and 9. Henceforth crafted by conveyance is sagaciously done by L2 norm distance vector estimation. Henceforth L2 is improved version of Euclidian for computing the centroid purposes of disseminated node bunches.

## VIII. CONCLUSION

In this paper we propose a novel mechanism for clustering biological information utilizing MPI and MapReduce with improved distance vector computation. The proposed mechanism adequately chooses the middle point when contrasted with existing mechanism and can convey information on various levels by making limit between them.

**REFERENCES**

[1] T. Estrada, B. Zhang, P. Cicotti, R. Armen, and M. Taufer, "Accurate analysis of large datasets of protein-ligand binding geometries using advanced clustering methods," Computers in Biology and Medicine, vol. 42, no. 7, pp. 758–771, 2012.

[2] J. Dean and S. Ghemawat, "MapReduce: simplified data processing on large clusters," Communications of the ACM, vol. 51, no. 1, pp. 107–113, 2008.

[3] [3] S. J. Plimpton and K. D. Devine, "MapReduce in MPI for large-scale graph algorithms," Parallel Computing, vol. 37, no. 9, Sep. 2011.

[4] A. Jain, "Bias, reporting, and sharing: Computational evaluations of docking methods," Journal of Computer-Aided Molecular Design, vol. 22, no. 3-4, pp. 201–212, 2008.

[5] S. S. Shende and A. D. Malony, "The TAU parallel performance system," International Journal of High Performance Computing Applications, vol. 20, no. 2, pp. 287–311, 2006.

[6] K. A. Abdul Nazeer, M. P. Sebastian,îImproving the Accuracy and Efficiency of thek-means Clustering Algorithm, Proceedings of the World Congress on Engineering 2009 Vol I WCE 2009, July 1 - 3, 2009, London, U.K.

[7] Soumi Ghosh, Sanjay Kumar Dubey, Comparative Analysis of K-Means and Fuzzy C-Means Algorithmsî, International Journal of Advanced Computer Science and Applications, Vol. 4, No.4, 2013

[8] Shafeeq,A., Hareesha,K.,ìDynamic Clustering of Data with Modified K-Means Algorithm, International Conference on Information and Computer Networks, vol. 27 ,2012

[9] Junatao Wang, XiaolongSu,îAn Improved K-means Clustering Algorithm, Communication Software and Networks (ICCSN), 2011 IEEE 3rd International Conference on 27 may,2011 (pp. 44-46)

[10] Shi Na, Liu Xumin, Guan Yong, ìResearch on K-means Clustering Algorithm: An Improved K-means Clustering Algorithm, Intelligent Information Technology and Security Informatics,2010 IEEE Third International Symposium on 2-4 April, 2010(pp. 63-67)

[11] H. L. Shashirekha and A. H. Wani, "Analysis of imputation algorithms for microarray gene expression data," 2015 International Conference on Applied and Theoretical Computing and Communication Technology (iCATccT), Davangere, 2015, pp. 589-593.

[12] A. H. Beg and M. Z. Islam, "Clustering by genetic algorithm- high quality chromosome selection for initial population," 2015 IEEE 10th Conference on Industrial Electronics and Applications (ICIEA), Auckland, 2015, pp. 129-134.