

# DynMR: A Dynamic Slot Allocation Framework for MapReduce Clusters in Big Data Management using DHSA and SEPB

Anil Sagar T<sup>1</sup>, Ramakrishna V Moni<sup>2</sup>

<sup>1</sup>(Mtech, Dept of CSE, VTU, SaIT, Bangalore

<sup>2</sup> (Prof, Dept of CSE, SaIT, Bangalore

\*\*\*\*\*

## Abstract:

MapReduce is one among the famous processing model for huge scale information (Big Data) processing in distributed computing. Since there may be a possibility of slot based MapReduce framework (eg. Hadoop MRv1) displaying some poor execution as a result of its unoptimized resource allocation. To venture on this, this paper finds and further streamlines the data distribution and resource allocation from the following three key perspectives. To begin with, because of the pre-configuration of the map slots and reduce slots which are not replaceable slots can be extremely under used. Since map slots may be completely used while reduce slots are empty and the other way around, considering the slot based model we set forth an option strategy called Dynamic Hadoop Slot Allocation. It unwinds the slot allocation parameters to permit slots to be reallocated to map or reduce task assignments relying upon their needs. Second the speculative execution can handle the straggler issue which sufficiently fit to enhance the execution for a job however to determine the expense of cluster proficiency. In context of this we further show Speculative Execution Performance Balancing so as to adjust the execution exchange between a single job and a batch of jobs. Third, delay scheduling has indicated to enhance the information and data locality at the fair cost. On the other hand we propose a method called Slot Pre Scheduling that can enhance the data locality yet with no effect on cost. At last by melding all the strategies together we make an orderly slot allocation framework called DynMR (Dynamic Map Reduce) which can enhance the execution of MapReduce workloads significantly.

**Keywords — MapReduce, DynMR, Delay Scheduler, Hadoop Fair Scheduler, Slot Allocation, Slot Pre Scheduler.**

\*\*\*\*\*

## I. INTRODUCTION

Despite the fact that, having numerous studies in improving MapReduce/Hadoop, there are couple of key difficulties for the usage and execution challenge of a Hadoop[4] cluster.

To answer these difficulties, we introduce DynMR, an element space allotment structure that enhances the execution of a MapReduce cluster through upgrading the opening use. Especially, Dynamic MR focuses over Hadoop Fair Scheduler

(HFS). This is on the grounds that the cluster usage and execution for MapReduce employments under HFS are much poorer (or more genuine) than that under FIFO scheduler. Actually side the indispensable purpose of our DynMR can be utilized for FIFO scheduler also. DynMR comprises of three optimization procedures, namely, Dynamic Hadoop Slot Allocation (DHSA), Speculative Execution Performance Balancing (SEPB) and Slot Pre-Scheduling from diverse key perspectives.

### ***Dynamic Hadoop Slot Allocation (DHSA)***

Just like YARN [10] which proposes another resource model of "container" that both map and reduce tasks can run on DHSA keeps the slot based resource model here DHSA will break the assumption of slot allocation as below

- slots are nonexclusive and even as it can be utilized by either map or reduce slots, although there is a pre-configuration for the quantity of map and reduced slots at the end of the day, when there are lacking map slots the map slots can use all the map slots and after that it can use the reduced slots available from the reduce slots also.
- Map slots prefer to utilize map slots, likewise reduce slots like to utilize reduce slots. The focal point is that, the pre-configuration of map and reduce slots for every slave node can still be used in any case to control the proportion of running map and reduced tasks during runtime, which can be better than YARN which has no control for the degree of running map and reduce tasks.

### ***Speculative Execution Performance Balancing (SEPB)***

Speculative execution is a basic procedure that can be utilized to address the issue of moderate running assignment's influence on a single job's execution time by running a reinforcement task and on another machine. In this paper, we propose a dynamic technique of slot allocation system called Speculative Execution Performance Balancing (SEPB) for the speculative job tasks. It can adjust the execution exchange off between a single job's execution time and a batch of jobs' execution time by deciding powerfully when it is time designates slots for speculative job tasks.

### ***Slot Pre-Scheduling***

To enhance the data information locality in MapReduce [9] delay scheduling has turned out to be a powerful approach. In the perspective of this, we propose an alternative procedure called named Slot Pre-Scheduling that has capacity to enhance

the data information locality, however it has no negative effect on reasonableness. It is accomplished to determine the burden of load balancing between slave nodes.

We have coordinated Dynamic MR into Hadoop (especially Apache Hadoop 1.2.1). We assess it utilizing proving ground workloads.

The fundamental commitments of this paper are abridged as follows:

- Propose a Dynamic Hadoop Slot Allocation (DHSA) strategy to increase the slot allocation for Hadoop.
- Propose a Speculative Execution Performance Balancing (SEPB) strategy to adjust the execution exchange off between a single jobs and a batch of jobs.
- Propose a Pre-Scheduling strategy to enhance the data information locality to the determination of load balance across the cluster, which has no negative influence on reasonableness.
- Develop a framework called DynMR by consolidating these three methods in Hadoop MRv1.
- Experiments have been performed to validate the effectiveness of Dynamic MR and its three step-by-step techniques.

## **II. RELATED WORK**

### ***Dynamic Split Model of Resource Utilization in MapReduce***

The prominence of MapReduce is expanding step by step as a parallel programming model for huge scale information preparing. At the same time in the long run, we discover some customary MapReduce stages which have a poor execution in content of bunch asset use following the conventional multi-stage parallel model and some current timetable approaches utilized as a part of the group environment have a few disadvantages. We address these issues through our involvement in planning a Dynamic Split Model of the assets usage which contains two advances, Dynamic Resource Allocation considering the stage need and employment prerequisite when distributing assets

and Resource Usage Pipeline which can relegate errands alertly.

The fundamental pipe line of parallel processing catches the scholastic world consideration. Further area/circle is an appropriated figuring stage which is like Google GFS/MapReduce. It comprises of a parallel runtime Sphere and in addition a conveyed file framework Sector. Another parallel runtime is phaser which is a direction build for element parallelism under the environment of multi-processors rather than multi-nodes.

The embodiment fundamental of these parallel programming edge works specified above is that they are all gotten from multi-stage parallel model. The customary multi-stage model has poor execution on the asset use efficiency. This issue has taken conception from two angles. From one viewpoint, different stages have different needs and also distinctive resources use inclination and must be executed entirely to that need which causes resource utilization unbalance. A few systems can release the strict execution request among distinctive stages by sub-operations cover execution, for example, phaser aggregator or gushing pipeline in Hadoop Online model. Anyways, the unbalance still exists if the resource require in sub-operations' cover execution is the same.

We came up with an innovation in this paper to address the above issue which is known as Dynamic Split Model of Resources Utilization which embodies two fundamental advances: Resource Usage Pipeline (RUP) and Dynamic Resource Allocation (DRA). What's more, to make the resource pipeline feasible we have to take a few measures to partition the mixture of different resources used to partitioned powerfully and dispatch an undertaking at a legitimate point. Another technique is that Dynamic Resource Allocation will number the framework load and the status of every job keeping in mind the end goal to partition our resource more efficiently. Also, we utilize Hadoop to confirm our advancement. Since our technique is executed on Hadoop, the issue we address is in no means constrained to Hadoop, even MapReduce processing model. By and large, effective resource use we can give a more significant execution in DISC system

### **Background**

One of the programming models which is appropriate in DISC system is called as MapReduce. The transitional yield of another sort of (key, quality) sets is created further and exchanged to reduce function. Reduce function will prepare all qualities fit in with the same key one time and yield the final (key, worth) sets.

Hadoop is a standout amongst the most well known open source executions of Google GFS/MapReduce. It is comprised of two sections Hadoop Distributed File System (HDFS) and MapReduce figure structure. MapReduce system is based on top of HDFS containing a Job Tracker and Task Tracker.

### **Map Task**

A map task relating to a specific occupation will read a part which is a segment of the input file from HDFS. The procedure is explained in the accompanying step

- MapRunner reads a (key, value) pair from the assigned split in HDFS using RecordReader and applies the user-defined map function to the pair.
- After processing one pair, map task outputs the result to a fixed-size buffer. If the buffer is overflow it will apply a quick-sort to the full buffer first and then store the content of the buffer to a spill file so the buffer becomes empty and can therefore receive more results. This process will take place repeatedly until all records are used up.
- Map task merges all the spills on the disk into a sort file and stores it to the local file system as well as a corresponding index file referring to that data file. Further, task tracker will ask for a new task from job tracker as soon as the map task finishes and makes a slot free.

### **Reduce Task**

A reduce task execution comprises three phases.

- In the shuffle phase reduce task fetches a specific partition of every map task output using HTTP request and puts them into memory first. If the memory is full it will apply merge sort to the memory and output the result to a temp file.

- After all partition is fetched, reduce task enters the sort phase where it sorts all sorted files stored in memory or disk using merge sort grouping all records to the same key together.
- The reduce phase applies the user-defined reduce function to each key and the set of values belonging to the same key.

### ***Dynamic Scheduling Model***

The real execution situation on a single node in the raw version Hadoop, Where the raw version Hadoop does have several flaws which throws a major impact on the system resource efficient utilization.

The genuine execution circumstance on a single node in the raw version of Hadoop, Where the raw variant Hadoop has a few flaws which tosses a real effect on the framework resource utilization

- Single node resource usage unbalance: inside a round of map task or reduce task different stages have different resource use predisposition particularly when homogenous undertakings execute at the same pace a few resource such as IO or CPU may be underused while the others abused.
- Reduce slot hoarding: In MapReduce, every Reduced tasks it partitions by the aftereffects of every map tasks, and can just apply the client's reduce tasks once it has results from all map tasks. Since, in the event that we present a difficult task, which will take quite a while to finish all map tasks. It will hold reduce slots for quite a while until all the map assignments are finished, and starve jobs occupations and underutilize resource.
- Resource allocation unbalance within job: MapReduce will assign resource by a static configuration, which does not consider the framework burden and the jobs necessity. This will prompt asset designation unbalance.
- We separate resource usage within a phase into two periods: CPU period and IO period. What's more, we utilize commercial advanced scheduling to dispatch a task at a legitimate point. One assignment's sub-operation can

covered with the other in the event that their resource use is complimentary.

- We gather the framework load and the status of every occupation at run time to dispense resource alertly. Further at the three subjective time point the quantity of slot is not the same and can be modified as per framework load.

### ***Dynamic Resource Allocation***

***Reduce Slot Hoarding Problem.*** MapReduce typically dispatches reduce tasks for a vocation when its first couple of maps finish, so these reduces can start replicating map yields while the remaining maps are running. In any case, in a substantial jobs, the map phase may take quite a while to finish. The employment will hold any reduce slots it gets amid this until its maps finish. So alternate occupations, which submitted later, will starve until the vast occupation finishes. This is called "reduce slot hoarding" issue, which will squander resource and deferral the divider time of jobs.

We may recognize that this issue can be explained by beginning Reduce tasks later or making them suspended, however called attention to that this arrangement is not achievable. And their answer is to part decrease undertakings into two intelligently particular sorts of tasks, duplicate assignments and process tasks, with partitioned types of control.

### ***Resources Allocation unbalance Problem***

Another issue we meet in MapReduce is that the MapReduce will distribute the most extreme number of slots by a static configuration, for example, the parameters: `mapred.tasktracker.maptasks.maximum` and `mapred.tasktracker.reduce.tasks`. most extreme, and will never be changed when the batch is running. Then again, the necessity for openings shifts alongside occupation continuing.

### ***Solution***

Dynamic Resource Allocation, Our proposed answer for these issues is dynamic resource allotment. We will assign resource as indicated by the group load and all occupations run-time status. As said before, amid the execution of an occupation, the resource prerequisite is changing without a

doubt. We oblige more map slots toward the starting; with the map tasks obliged map tasks is bit by bit reduced, yet the prerequisite to reduce slots expanded; And more Reduced slots are required when all map tasks are finished. Thus, for a vocation, the resource prerequisite is changing with occupation status evolving. So we define the element weight of the guide stage and the reduce phase as indicated by the occupation status to recreate the element prerequisite for resources.

### III. OVER VIEW

Further, to enhance the execution of a MapReduce cluster through optimizing the slot usage fundamentally from two points of view. In the first place, there are two separate slots, to be specific, busy slots (i.e., with running tasks) and idle slots (i.e., no running tasks). Given the aggregate number of map and reduce slots configured by clients, one optimizing methodology (i.e., macro-level advancement) is to enhance the slot used by amplifying the quantity of occupied slots and decreasing the quantity of idle slots. Second, it is important that not every occupied opening can be efficiently used. Consequently, our optimization approach (i.e., micro-level enhancement) coordinates towards the effecient use of busy slots occupied after the macro- level optimization. We propose Dynamic MR, an element usage optimization structure for MapReduce, to enhance the execution of an imparted Hadoop cluster under a fair scheduling between clients

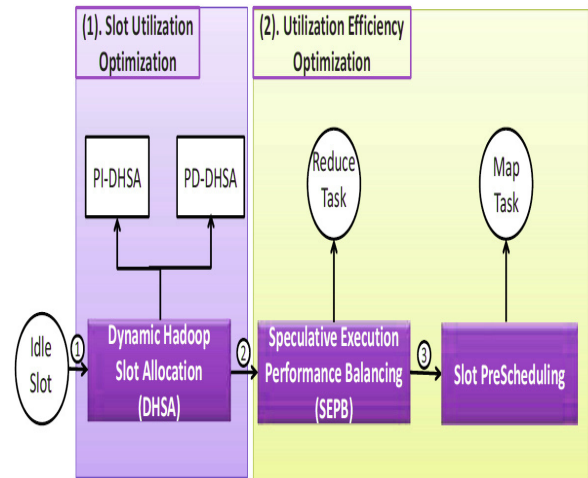


Fig. 1: DynamicMR Framework.

Figure 1 gives an overview of DynamicMR. It consists of three slot allocation techniques, i.e., *Dynamic Hadoop Slot Allocation (DHSA)*, *Speculative Execution Performance Balancing (SEPB)*, and *Slot PreScheduling*.

Every single system considers the execution performance from diverse perspectives. DHSA endeavor to minimize slots usage while keeping up the fairness, when there are pending tasks (e.g., map tasks or reduced tasks). SEPB recognizes the slot resource in-efficiency issue for a Hadoop cluster, brought about by speculative tasks. slot Pre-Scheduling enhances the slot utilization efficiency and execution performance by enhancing the data locality for map tasks while keeping the fairness.

By fusing the three systems, it empowers DynamicMR to upgrade the use and execution of a Hadoop cluster significantly with the accompanying regulated procedures

- Whenever there is an idel slot accessible, DynamicMR will first try to enhance the slot used with DHSA. It chooses alertly whether to allocate it or not, subject to the various constraints, e.g., fairness and load balancing.
- If the distribution is true, Dynamic MR will further enhance the execution by enhancing the effectiveness of slot utilization with SEPB. Since the speculative execution can enhance the

execution of a single job however to the detriment of cluster effectiveness, SEPB goes about as a productivity balance between a single job and a cluster efficiency. It takes a shot at top of Hadoop speculative scheduler to focus rapidly whether allocating the idle slots to the pending tasks or speculative tasks.

- When allocating the idle slots for pending/speculative map tasks, Dynamic MR will have the capacity to further enhance the slot utilization efficiency from the data locality optimization aspects with Slot Pre-Scheduling.

Besides, we need to specify that the three procedures are at diverse levels, i.e., they can be connected together or independently.

DynamicMR components can be explained in detail in the further discussions.

**Dynamic Hadoop Slot Allocation (DHSA)**

The current configuration of MapReduce experiences an under-usage of the slots as the quantity of map and reduce tasks shifts over the long run. Our dynamic slot allocation approach is taking into account the perception that at distinctive time there may be idle map(or reduce) slots, as the jobs continues from map stage to reduce stage. We can utilize the unused map slots for those overburden reduce tasks to enhance the execution of the MapReduce workload, and the other way around. We further make utilization of idle reduce slots for running map tasks. That is, we break the certain presumption for current MapReduce structure that the map tasks can just run on map slots and reduced tasks can just run on reduce slots.

There are two challenges specified below that must be considered:

(C1): Fairness is an imperative metric in Hadoop Fair Scheduler (HFS). We proclaim it as reasonable when all pools have been designated with the same amount of resource. In HFS, task slots are first allocated over the pools [8], and later then the slots are distributed to the jobs inside the pool. Also, a MapReduce job computation embodies two sections: map-phase task computation and reduce-phase task computation.

(C2): The resource requirement between the map slots and reduced slots are especially diverse. The purpose for this is the map tasks and reduced tasks

regularly show totally different execution designs. reduce task has a tendency to expend considerably more resources, for example, memory and system network speed. Basically permitting reduce tasks to utilize map slots configuring every map slots to take more resources, which will therefore lessen the powerful number of slots on every node, creating resources under-used amid runtime.

With a due appreciation towards (C1), we set forth a Dynamic Hadoop Slot Allocation (DHSA). It contains two choices, to be specific, pool- free DHSA(PI-DHSA)

**pool-Independent DHSA (PI-DHSA)**

HFS utilizes max-min fairness [5] to allocate slots crosswise over pools with least ensures at the map-phase and reduce-phase, individually. Pool-Independent DHSA (PI-DHSA) extends the HFS by dispensing slots from the clusters of worldwide level and free of pools.

The allocation procedure is comprised of two sections:

- Intra-phase dynamic slot allocation: Each pool is part into two sub-pools, i.e., map phase pool and reduce phase pool. At every stage, every pool will get its share of slots.
- Inter-phase dynamic slot allocation: After the intra-phase dynamic slot allocation for both the map-phase and reduced phase, next we can perform the dynamic slot allocation crosswise over typed phase.

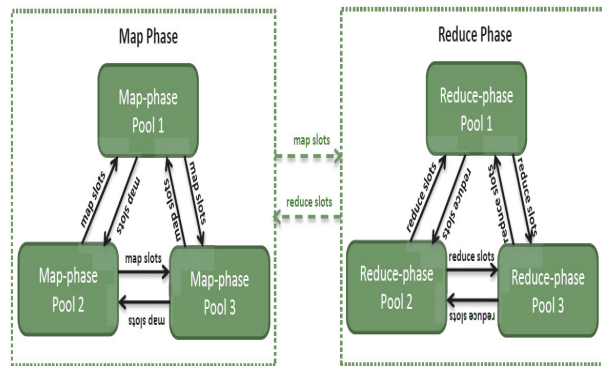


Fig. 2: fairness-based slot allocation flow for PIDHSA.

The entire dynamic slot allocation flow is that, at whatever point a pulse is gotten from a computing node, at first we process the aggregate demand for map slots and reduce slots for the current MapReduce workload. At that point we focus alertly the need to acquire map (or reduce) slots for reduce (or map) tasks in light of the interest for map and reduce slots, with respect to these four situations. The specific number of map (or reduce) slots to be obtained is based on the account of quantity of unused reduced (or map) slots and its map (or reduce) slots needed.

To accomplish the reservation usefulness, we give two variables rate Of Borrowed Map Slots and rate Of- Borrowed Reduce Slots, defined as the rate of unused map and reduced slots that can be obtained, separately. Thus, we can restrict the quantity of unused map and reduced slots that ought to be distributed for map and reduced tasks at every pulse of that task tracker. With these two parameters, clients can flexibly adjust the exchange off between the performance execution optimization and the starvation minimization.

In addition, Challenge (C2) makes us to review that we can't treat map and reduce slots as same, and just obtain unused slots for map and reduce tasks. Rather, we should be mindful of shifted resource sizes of map and reduce slots. A slot weight- based methodology is therefore proposed to address the issue. We allot the map and reduce slots with distinctive weight values, regarding the asset configurations. Particular to the weights, we can alterably decide the amount of map and reduce tasks which has to be generate in the length of runtime.

**Pool-Dependent DHSA (PD-DHSA)**

As an opposite point on checking towards PI-DHSA Pool-Dependent DHSA (PD-DHSA) considers fairness for the dynamic slot allocation across pools. Accepting that every pool, includes two sections: Map phase pool and Dynamic Phase pool, is selfish. It is considered fair when aggregate quantities of map and reduce slots allocated across pools are the same with one another. PD-DHSA will be performed with the accompanying two courses of actions:

(1). *Intra-pool dynamic slot allocation.* At a early stage, each typed- phase pool will receive its share of typed-slots based on max-min fairness at each phase. There are four possible relationships cases for every pool regarding its demand (denoted as mapSlots Demand, reduceSlots Demand) and its workload (marked as mapShare, reduceShare) between two phases:

Case (a).  $mapSlotsDemand < reduceShare$ , and  $reduceSlots-Demand > reduceShare$ . We can use some of the unused map slots for its overloaded reduce tasks from its reduce-phase pool first before using other pools.

Case (b).  $mapSlotsDemand > mapShare$ , and  $reduceSlots- Demand < reduceShare$ . we can use some unused reduce slots for its map tasks from its map-phase pool first before using pools.

Case (c).  $mapSlotsDemand < mapShare$ , and  $reduceSlots- Demand < reduceShare$ . Both map slots and reduce slots are enough for its use. It can give some unused map slots and reduce slots to other pools.

Case (d).  $mapSlotsDemand > mapShare$ , and  $reduceSlots- Demand > reduceShare$ . If both map slots and reduce slots of a pool have become insufficient. It may have to borrow some unused map or reduce slots from other pools through inter-Pool dynamic slot allocation is shown below.

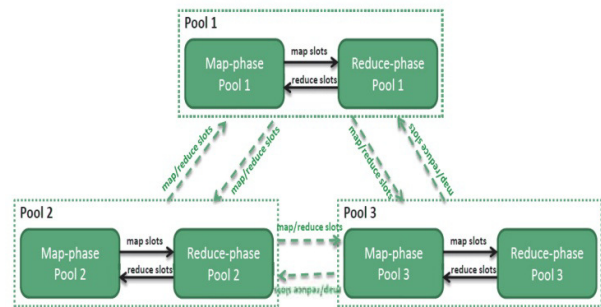


Fig. 3: Example of the fairness-based slot allocation flow for PD-DHSA. The black arrow line and dash line show the borrow flow for slots across pools.

(2). *Inter-pool dynamic slot allocation.* It is obvious that,

(i). if a pool, has  $mapSlotsDemand + reduceSlotsDemand < mapShare + reduceShare$ . The slots are enough for the pool and there is no

need to get some map or reduce slots from other pools

(ii). On the contrary, when  $\text{mapSlotsDemand} + \text{reduceSlotsDemand} \leq \text{mapShare} + \text{reduceShare}$ , the slots are not enough even after Intra-pool dynamic slot allocation.

The overall slot allocation process for PD-DHSA is as sketched down below in figure 4.

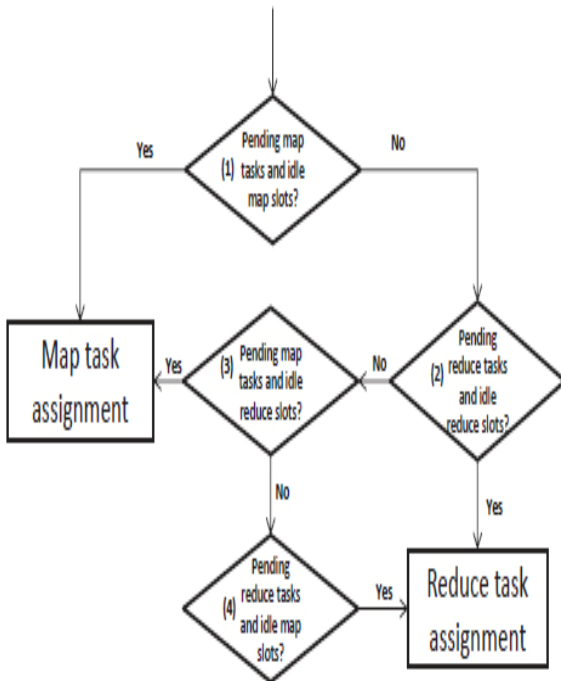


Fig. 4: The slot allocation flow for each pool under PD-DHSA. The numbers labeled in the graph corresponds to Case (1)-(4) in Section 2.1.2, respectively.

At first, it computes the maximum number of free slots that can be allocated at each round of heartbeat for the tasktracker. Next it starts the slot allocation for pools. For every pool, there are four possible slot allocations as illustrated in Figure 4 above.

Case(1): We try the map tasks allocation, if there are idle map slots for the task tracker, and there are pending map tasks for the pool.

Case(2): If the attempt of Case(1) fails, the condition does not hold good, and it cannot find a map task satisfying the valid data-locality level, we continue to try reduce tasks allocation when there are pending reduce tasks and idle reduce slots.

Case(3): If Case(2) fails due to the required conditions does not hold, we try for map task allocation again. If Case(1) fails then there might not have to be any idle map slots available. In contrast, if Case(2) fails then there are no pending reduce tasks. In this case, we can relay on reduce slots for map tasks of the pool.

Case(4): If Case(3) fails, we try for reduce task allocation once again. Case(1) and Case(3) fail might be because of no valid locality-level pending and map tasks available, but there are idle map slots. In contrast, Case(2) might not have any idle reduce slots available. At such cases, we can allocate map slots for reduce tasks for the pool.

Furthermore, there is a special scenario that needs to be considered particularly. Note, it is possible that all the above four possible slot allocation attempts fail for all pools, due to the data locality consideration for map tasks.

### Speculative Execution Performance Balancing (SEPB)

MapReduce job’s execution time is very sensitive to slow- running tasks (namely straggler) [6], [7]. We divide the stragglers into two types, namely Hard Straggler and Soft Straggler, as defined below:

- **Hard Straggler:** A task that gets into deadlock status due to the endless waiting for some resources. It cannot stop and complete unless we stop it manually.
- **Soft Straggler:** A task that can complete its task successfully, but will take much longer time than that of common tasks.

For the hard straggler, we should stop it and run another task, or called a *backup task*, immediately once it was detected. In contrast, there are two possible cases between the soft straggler and its backup task:

(S1). Soft straggler finishes first or at the same time as its backup task. For this case, there is no need to run backup task.

(S2). If Soft straggler finishes later than the backup task. We should stop it and run a backup task immediately.

Further to deal with the straggler problem, speculative execution is used in Hadoop. Instead of diagnosing and fixing straggling tasks, it finds the straggling task dynamically using heuristic algorithms such as LATE [7]. Once detected,



however, it can- not simply kill the straggler immediately due to the following facts:

- Hadoop does not have a way or methodology to distinguish between the hard straggler and the soft straggler.
- But for the soft straggler, it's also difficult to judge whether it belongs to (S1) or (S2) before running. Simply stopping the straggler will harm the case of (S1).

Rather, it produces a backup task and permits it to run simultaneously with the straggler, i.e., there is a processing cover between the straggler and the backup undertaking. The task killing operation will happen when both of the two tasks are finished. It merits specifying that, despite the fact that the speculative execution can reduce a single work's execution time, however it has a go at the expense of cluster effectiveness.

Along these lines, it raises a test issue for speculative tasks on the best way to relieve its negative effect for the execution of batch jobs. To expand the execution for a group of jobs, an instinctive arrangement is that, given accessible task slots, we should fulfil pending tasks first before considering speculative tasks. That is, the point at which a node has an idle map slots , we have to pick pending map tasks first before searching for speculative map tasks for a cluster of jobs.

We further propose a dynamic task allocation mechanism called Speculative Execution Performance Balancing (SEPB) for a batch of jobs with speculative execution tasks on top of Hadoop's current task selection strategy. Hadoop picks a task from a job in view of the accompanying need: first, any failed task is given the most highest priority. Second, the pending tasks are considered. For map, tasks with data local to the process node are picked first. At last, Hadoop searches for a straggling assignment to execute speculatively. In our task scheduling component, we define a variable rate Of Jobs Checked- For Pending Tasks with domain somewhere around 0.0 and 1.0, configurable by clients, to control max Num Of Jobs Checked For- Pending Tasks, which is the greatest number of occupations that are checked for pending map and reduced taskd for a batch of jobs.

Then again, we can perceive another challenging issue if there is a delay in the planning of speculative task. To defeat this issue, at present we utilize a basic heuristic algorithm: We evaluate the execution time for every task. When it took twice more than the normal execution time of tasks, we kill it specifically to yield the slot. Since failed/executed tasks have the most elevated need to run in Hadoop, a reinforcement undertaking will be made to supplant it rapidly, enhancing the execution of a single job and moderating the negative effect on the cluster effectiveness.

### ***Discussion on SEPB VS LATE***

The benefit of SEPB over LATE lies in its arrangement for slot allocation to speculative tasks. Conversely, SEPB performs the resource allocation for speculative tasks from a worldwide view by considering various occupations (controlled by the argument max Num of Jobs Checked for Pending Tasks). Further postpones the slot allocation to speculative tasks at whatever point there are pending tasks for the different jobs. The SEPB figures out if to make a speculative task to re-figure information or not from a global view by checking various jobs. On the off chance that SEPB recognizes pending tasks, it will assign the idel slots to a pending tasks. If not, another speculative task will then be made to have the idle slots.

### ***Slot PreScheduling***

We propose a Slot Pre-Scheduling technique that holds ability to improve the data locality while having no negative impact on the fairness of MapReduce jobs. The basic level idea is that, in light of the fact that there are often some idle slots which cannot be allocated due to the load balancing constrain during runtime, we can pre-allocate those slots of the node to jobs to maximize the data locality.

We propose a Slot Pre-Scheduling system that holds capacity to enhance the data locality while having no negative effect on the fairness of MapReduce jobs. The essential level idea is that, in a way there are regularly some idle slots which can't be dispensed because of the load balancing during runtime, we can preallocate those slots of the node to allocate and amplify the data locality

### **Preliminary**

Prior to presenting Slot PreScheduling, we start with two definitions:

**Definition 1.** The **allowable idle map (or reduce) slots** will relate to the maximum number of idle map (or reduce) slots that can be allocated for a task tracker, considering the load balancing between machines.

**Definition 2.** The **extra idle map (or reduce) slots** will relate to the remaining idle map (or reduce) slots by subtracting the maximum value of used map (or reduce) slots and allowable idle map (or reduce) slots from the total number of map slots for a task tracker, considering the load balancing between machines.

### **Observation and Optimization**

For a MapReduce cluster, the computing workloads of running map (or reduce) tasks between task trackers (i.e., machines) are generally modified, because of the following facts.

(1) In practical world, Lots of MapReduce clusters comprises of heterogeneous machines (i.e., different computing powers between machines).

(2) There are often varied computing loads (i.e., execution time) for map and reduce tasks from different jobs, due to the varied input data sizes as well as applications.

(3) Even for a single job under the homogenous environment, the execution time for map (or reduce) tasks will not be the same.

In order to balance the workload, Hadoop comes up with a methodology that can dynamically control the number of allowable idle map (or reduce) slots (See **Definition 1**) for a task tracker in a heartbeat as the following three steps.

**Step 1#:** Compute the load factor  $mapSlotsLoadFactor$  as the sum of pending map tasks and running map tasks from all jobs divided by the cluster map slot capacity.

**Step 2#:** Compute the current maximum number of usable map slots by multiplying  $\min\{mapSlotsLoadFactor, 1\}$  with the number of map slots in a task tracker.

**Step 3#:** Finally, we can compute the current allowable idle map (or reduce) slots for a task tracker, by subtracting the current number of used

map (or reduce) slots from the current maximum number of usable map slots.

To make use of Slot Pre-Scheduling there are two different cases. The first case considers a task tracker slot on which there are extra idle map slots available, but no allowable idle map slots. For a headed job following the fair-share priority order, when it has local map tasks with block data on the task tracker slot, instead of skipping it by the default Hadoop scheduler, we can proactively allocate extra map slots to the job.

The second case is for DHSA. When there are no idle map slots but some idle reduce slots available, for a connected task tracker slot in a heartbeat, we can proactively borrow idle reduce slots for local pending map tasks and restore them later, in order to maximize the data locality.

### **Advantages**

- Improves the performance of MapReduce workloads write maintaining the fairness.
- Can be used for any kinds of MapReduce jobs (independent or dependent ones).
- Balances the performance trade-off between a single job & a batch of jobs dynamically.
- Slot pre-scheduling improves the efficiency of slot utilization by further maximizing its data locality.
- Dynamic MR improves the performance of the Hadoop system significantly.
- SEPB identify the slot inefficiency problem of speculative execution.
- Dynamic MR consistently outperforms YARN.

### **Disadvantages**

- On comparison with YARN, the experiments show that, for single jobs, the result is inconclusive.
- The proposed Dynamic MR doesn't considers the implementation on cloud computing environment which is a gateway for further research.

#### IV. RESULTS AND ANALYSIS

##### Experimental Set up

We ran our experiments in a group comprising of 10 process hubs, each with two Intel CPUs (4 CPU centres every CPU with 3.07 GHz), 24GB memory and 56GB hard disk. We arrange one node as master and namenode, and the other 9 nodes as slaves and datanodes. The most recent rendition of Hadoop 1.2.1 is picked in our experiment.

##### Performance improvement for DynMR

In this segment, we first demonstrate the execution forms for PI-DHSA and PD-DHSA. At that point we assess and look at the execution change by PI-DHSA and PDDHSA under distinctive slot setup. Third, we make a dialog on the execution impact of the contentions of the rate of map and reduce slots that can be obtained for our DHSA in Appendix F of the supplemental material.

To show distinctive levels of fairness for the dynamic task allocation calculations, PI-DHSA and PD-DHSA, we perform an analysis by considering three pools, each with one task submitted. Figure 5 demonstrates the execution stream for the two DHSAs, with 10 sec every time step. The quantity of running map and reduce task for every pool at every time step is recorded. For PI-DHSA, as delineated in Figure 5(a), we can see that, toward the starting, there are just map tasks, with all slots utilized by map slots under PI-DHSA.

Every pool imparts 1/3 of the aggregate openings (i.e., 36 spaces out of 108 openings), until the 5th time step. The map slots interest for Pool 1 starts to therapist and the unused map slots of its impart are respected Pool 2 and Pool 3 from the 6th time step. Next from 6th to 10th time step, the guide assignments from Pool 2 and Pool 3 just as impart all guide openings and the reduced tasks from Pool 1 have all reduced tasks, taking into account the write stage level fairness approach of PI-DHSA(i.e., intraphase element space distribution). From 11th to 18th time venture, there are some unused map slots from Pool 2 and they are controlled by map tasks from Pool 3 (i.e., intra-stage element opening portion). Later, there are some unused map slots from Pool 3 and they are utilized by reduced tasks from Pool 1 and Pool 2 from 22st to 25th time step (i.e., between stage dynamic space portion). For PD-DHSA, like PI-DHSA toward the starting, every pool gets 1/3 of the aggregate spaces from the 1th to 5rd time venture, as indicated in Figure 5(b). Some unused map slots from Pool 1 are respected Pool 2 and Pool 3 from 6th to the 7th time step. Be that as it may, from the 8th to 12th, the map tasks from Pool 2 and Pool 3 and the decrease assignments from Pool 1 takes 1/3 of the total slots, subject to the pool-level reasonableness strategy of PD-DHSA (i.e., intra-pool element space allotment). At long last, the unused slots from Pool 1 starts to respect Pool 2 and Pool 3 since 13th time step (i.e., between pool element space portion).

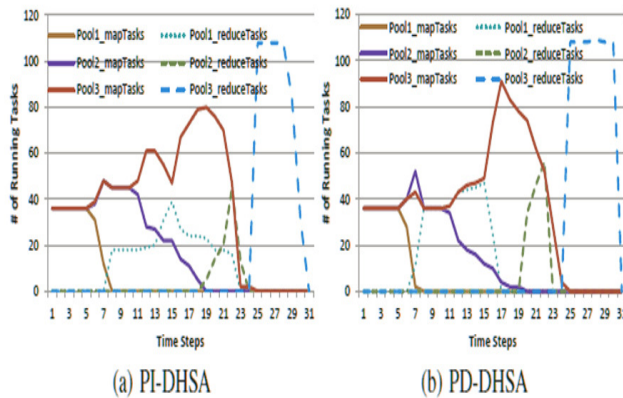


Fig. 5: The execution flow for the two DHSAs. There are three pools, with one running job each.

##### Speculative Execution Control for Performance

we expressed that theoretical assignment execution can defeat the issue of straggler (i.e., the moderate running errand) for an occupation, however it is at the expense of cluster use. We characterize a client's configurable variable. `percentageOfJobsCheckedForPendingTasks` to focus the time to calendar speculative tasks. To accept the adequacy of our element speculative execution control arrangement, we perform an investigation with 5 employments, 10 occupations and 20 occupations by fluctuating the estimations of `percentageOfJobsCheckedForPendingTasks`.

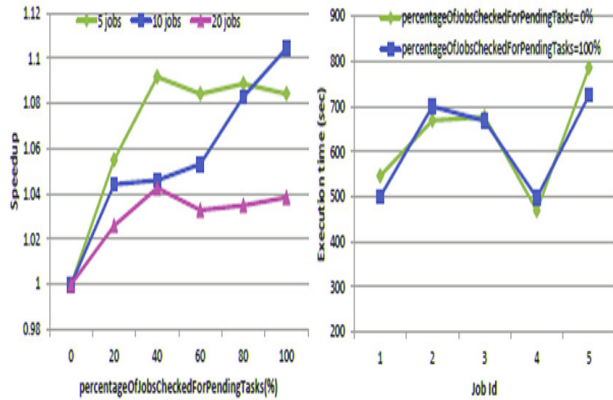


Fig. 6: The performance results with SEPB.

Note that LATE [35] has been actualized in 1.2.1. Figure 6 give the execution results SEPB in examination to LATE. All speedups are processed regarding the case that percentageOfJobsCheckedForPendingTasks is equivalent to zero.

We have the accompanying discoveries:

First and foremost, SEPB can enhance the execution of Hadoop from 3%-10%, indicated in Figure 6(a). As the estimation of percentageOfJobsCheckedForPendingTasks expands, the pattern of execution change has a tendency to be huge and the ideal setups could be unmistakable for diverse workloads. For instance, the ideal setup for 5 occupations is 80%, yet for 10 employments is 100%. The reason is that, huge estimation of rate-OfJobsCheckedForPendingTasks will let more quantities of employments be checked for pending assignments before considering speculative execution for every space assignment, i.e., It is more inclined to allot an opening to a pending tasks first and foremost, instead of a theoretical tasks, which advantages more for the entire jobs. Notwithstanding, extensive estimation of percentageOfJobsCheckedForPendingTasks will defer the speculative execution for straggled jobs, harming their execution. For a few workloads, too vast estimation of percentageOfJobsCheckedForPendingTasks will corrupt the execution for straggled jobs a great deal and thus influence the general occupations, clarifying why the ideal design is not generally

100%. We prescribe clients to design percentageOfJobsCheckedForPendingTasks at 60%-100% for their workloads.

Second, there is an execution trade-off between an individual job and the entire jobs with SEPB. We demonstrate a case for the workload of 5 occupations when setting percentageOfJobsChecked- ForPendingTasks to be 0 and 100%, separately. As results indicated in Figure 6(b), Job 2 and 4 are negative influenced because of the oblige on speculative execution from SEPB, though it supports the execution for entire jobs (i.e., the most extreme execution time of jobs)

### Data Locality Improvement Evaluation for Slot PreScheduling

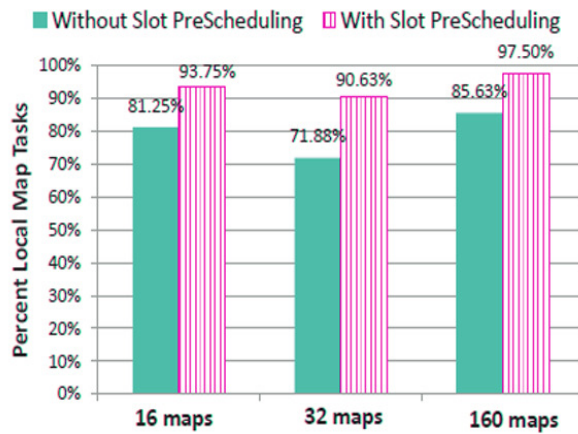


Fig. 7: The data locality improvement by Slot PreScheduling for Sort benchmark.

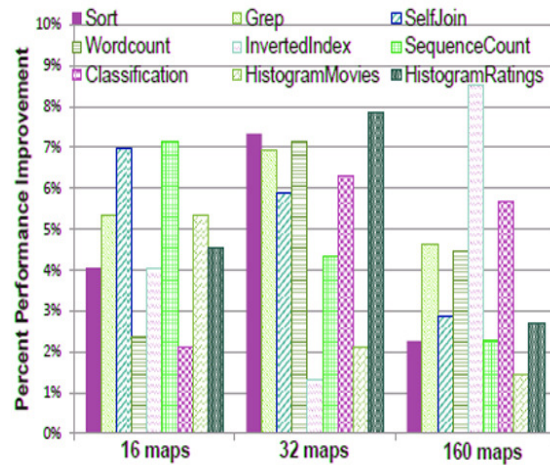


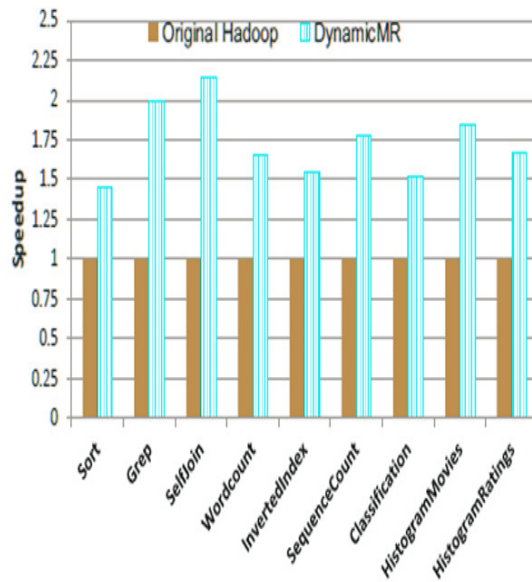
Fig. 8: The performance improvement under Slot PreScheduling.

To test the impact of Slot PreScheduling on data locality improvement, we ran MapReduce employments with 16, 32, and 160 map tasks on the Hadoop cluster. We contrast reasonable imparting results and without Slot PreScheduling under the default HFS. It merits saying that Delay Scheduler has been added to the default HFS for the customary Hadoop and continues working dependably. Subsequently, our work turns to be the correlation between the case with Delay Scheduler just and the case with Delay Scheduler in addition to Slot PreScheduling. Figure 7 demonstrates the data locality results with and without Slot PreScheduling for Sort benchmark. With Slot PreScheduling, there are around 2% ~25% region change on top of Delay Scheduler for Sort benchmark. Figure 8 exhibits the comparing execution results profiting from the information area change made by Slot PreScheduling. There are around 1% ~ 9% execution change concerning the first Hadoop for the previously stated 9 benchmarks separately. Also, we measure and analyse the heap uneven degree and shameful degree for Hadoop cluster with and without Slot PreScheduling in Appendix E of the supplemental material

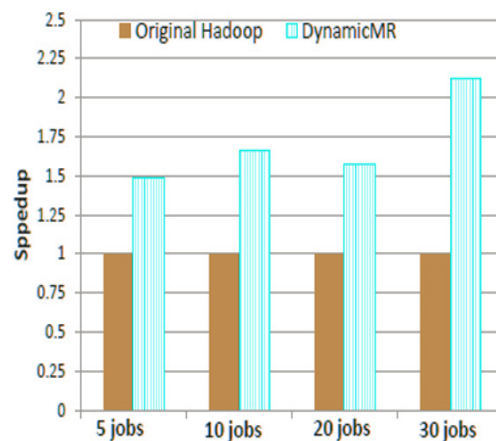
**Performance Improvement for DynamicMR**

In this segment, we assess DynamicMR framework in general by empowering all its three sub-schedulers with the goal that they can work corporately to augment the execution however much as could reasonably be expected. For DHSA part, we subjectively pick PI-DHSA, taking note of that PI-DHSA and PD-DHSA have fundamentally the same execution change (See Section 3.2.2). For the first Hadoop, we pick the ideal space design for MapReduce occupations by counting all the conceivable slot setups. We intend to contrast the execution for DynamicMR and the first Hadoop under the ideal map/reduce slot design for MapReduce jobs. Figure 9 exhibits the assessment results for a single MapReduce work and additionally MapReduce workloads comprising of numerous jobs. Especially, for different jobs, we consider 5 jobs, 10 jobs, 20 jobs, and 30 jobs under a clump accommodation, i.e., all jobs submitted in the meantime. All speedups are computed as for the

first Hadoop. We can see that, even under the advanced map/reduce slot arrangement for the first Hadoop, our DynamicMR framework can at present further enhance the execution of MapReduce jobs altogether, i.e., there are around 46% ~ 115% for a single jobs and 49% ~ 112% for MapReduce workloads with numerous jobs. Also, we likewise actualize our DynamicMR for Hadoop FIFO scheduler. To approve the adequacy of our DynamicMR, we perform explores different avenues regarding the previously stated MapReduce workloads.



(a) A single MapReduce job



(b) MapReduce workloads with multiple jobs

Fig. 9: The performance improvement with our DynamicMR system for MapReduce workloads.

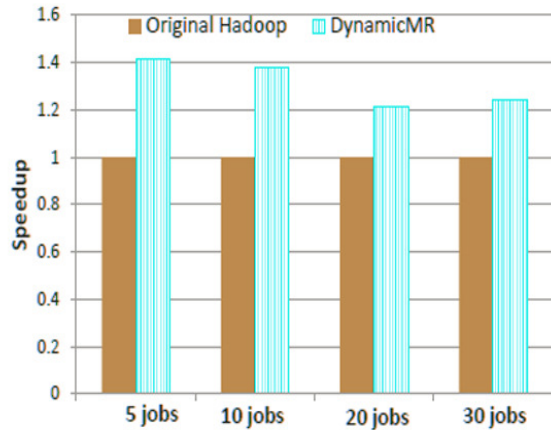


Fig. 10: The performance improvement with our DynamicMR system for MapReduce workloads under Hadoop FIFO scheduler.

## V. CONCLUSION

This paper proposes a DynamicMR Technique can be used to enhance the execution of MapReduce workloads while keeping up the fairness. It comprises of three methods, in particular, DHSA, SEPB, Slot PreScheduling, all of which concentrate on the slot use optimization for MapReduce group from alternate points of view. DHSA concentrates on the slot use expansion by distributing map or reduce slots to map and reduce tasks alterably. Especially, it doesn't have any presumption or require any earlier learning and can be utilized for any sorts of MapReduce jobs (e.g., autonomous or subordinate ones). Two sorts of DHSA are introduced, in particular, PI-DHSA and PD-DHSA, in view of distinctive levels of fairness. Client can pick both of them likewise. Rather than DHSA, SEPB and Slot PreScheduling consider the effectiveness advancement for a given slot usage. SEPB recognizes the slot unused issue of speculative execution. It can adjust the execution tradeoff between a single job and a batch of job alterably. Slot PreScheduling enhances the proficiency of slot use by expanding its data locality. By empowering the over three systems to work helpfully, the exploratory results demonstrate that our proposed DynamicMR can enhance the execution of the Hadoop framework altogether (i.e., 46% ~ 115% for single occupations and 49% ~ 112%

for various employments). In future, we plan to consider executing DynamicMR for distributed computing environment with more measurements (e.g., plan, due date) considered and distinctive stages by assessing some current works, [7], [10],[6].

## REFERENCES

- [1] F. Ahmad, S. Y. Lee, M. Thottethodi+, T. N. Vijaykumar. *PUMA: Purdue MapReduce Benchmarks Suite*. ECE Technical Reports, 2012.
- [2] G. Ananthanarayanan, S. Kandula, A. Greenberg, I. Stoica, Y. Lu, B. Saha, and E. Harris, *Reining in the outliers in map-reduce clusters using mantri*, in OSDI'10, pp. 1-16, 2010.
- [3] J. Dean and S. Ghemawat. *MapReduce: Simplified Data Processing on Large Clusters*, In OSDI'04, pp. 107-113, 2004.
- [4] Hadoop. <http://hadoop.apache.org>.
- [5] Max-Min Fairness (Wikipedia). [http://en.wikipedia.org/wiki/Max-min\\_fairness](http://en.wikipedia.org/wiki/Max-min_fairness).
- [6] T. White. *Hadoop: The Definitive Guide, 3rd Version*. O'Reilly Media, 2012.
- [7] M. Zaharia, A. Konwinski, A.D. Joseph, R. Katz, I. Stoica, *Improving MapReduce performance in heterogeneous environments*. In OSDI'08, pp.29-42, 2008.
- [8] M. Zaharia, D. Borthakur, J. Sarma, K. Elmeleegy, S. Schenker, I. Stoica, *Job Scheduling for Multi-user Mapreduce Clusters*. Technical Report EECS-2009-55, UC Berkeley Technical Report (2009).
- [9] M. Zaharia, D. Borthakur, J. Sarma, K. Elmeleegy, S. Schenker, I. Stoica, *Delay scheduling: A simple technique for achieving locality and fairness in cluster scheduling*. In EuroSys'10, pp. 265-278, 2010.
- [10] Shanjiang Tang, Bu-Sung Lee, Bingsheng He, *DynamicMR: A Dynamic Slot Allocation Optimization Framework for MapReduce Clusters*, IEEE Transactions On Cloud Computing



**Anil Sagar T\*** received his BE degree in Computer science and Engineering from Visveswaraya Technological University, Belgam, in 2013. Currently working towards the MTech degree in computer science at Sambhram Institute of Technology, Bangalore. His research interests are Big Data and Cloud Computing.



**Ramakrishna V Moni\*\*** currently working as prof in Department of CSE, Sambhram Institute of Technology, Bangalore. He received his B.Tech degree in 1986, M.S in 2000, Aero EnggC and PhD degree in the year 2007. His research interests are Big Data and Cloud Computing.