

# Augmenting Amdahl's Second Law: A Theoretical Model for Cost-Effective Balanced HPC Infrastructure for Data-Driven Science

Arghya Kusum Das, Jaeki Hong, Kisung Lee, Seung-Jong Park, Wooseok Chang

Department of Computer Science and Electrical Engineering, Center for Computation and Technology, Louisiana State University, Baton Rouge, LA, USA, 70803

Samsung Electronics Co., Ltd., 95, Samsung 2-ro, Giheung-gu, Yongin-si, Gyeonggi-do, S. Korea, 446711

**Abstract**—Current compute technologies to analyze enterprise and scientific big data started demanding more processors, storage, and memory resources forcing a similar growth in the total computation cost. Consequently, the HPC system designers are facing relentless pressure to reduce the system cost while providing expected performance also for these data- and compute-intensive applications. With this motivation, this paper proposes a simple, easy to use, additive model to optimize cost/performance by quantifying the system balance among CPU, I/O bandwidth and size of DRAM in terms of both application characteristics and the hardware cost. Our performance model theoretically augment Amdahl's second law for balanced system (Amdahl's I/O and memory number) and reveals that a balanced system needs almost 0.17GBPS I/O bandwidth, and almost 3GB of DRAM per GHz of CPU speed considering Intel Xeon processing architecture and current price trend in processor, storage and memory.

To substantiate our claim, we evaluate three fundamentally different cluster architecture, 1) a traditional HPC cluster called SupermikeII, 2) a regular datacenter called SwatIII, and 3) a novel micro brick based hyperscale system called CeresII. CeresII has 6-Xeon cores each running at 2GHz, 1-NVMe SSD with 2GBPS I/O bandwidth and 64GB DRAM, thus closely resembles to the optimum produced by our model. We evaluated these three clusters with respect to two widely used Hadoop-benchmark (terasort and wordcount) as well as our own genome assembler based on Hadoop and Giraph which serves as a real world example of data- and compute-intensive workload. CeresII outperformed all other cluster architecture for all the benchmark in terms of both execution time and cost/performance. For a large human genome assembly CeresII showed more than 85% improvement over SuperMikeII and almost 40% improvement over SwatIII in terms of cost/performance.

**Index Terms**—Cost to Performance, Balanced System, Hadoop, Giraph.

## I. INTRODUCTION

CURRENT compute technologies for enterprise and scientific big data analysis are demanding more compute cycles per processor, with extreme I/O performance also required. At the same time, the hardware cost for storing and processing these big data is increasing linearly with the increase in volume and complexity of these data. Consequently, today's system designers are facing relentless pressure

to reduce the system cost while providing the expected level of performance also. Complicating the scenario, the exponential growth of big data is rapidly changing the fundamental model of computation involved in high performance computing. A growing number of codes in both enterprise and scientific domain are being written using state of the art big data analytic software such as, Hadoop, Giraph which carefully consider data locality.

As a consequence, most of the existing high performance computing (HPC) clusters focusing only on tera to peta FLOP scale processing speed are found to be unbalanced ([1], [2]) for data-intensive scientific applications either in terms of performance, or economy or both. The existing data center architectures also needs to be reexamined. Although the current state of the art big data analytic software were initially developed to work well on scaled out cluster of cheap commodity hardware, recently published results (e.g., [3]) showed that scaled up can also provide better result for different data-intensive applications both in terms of performance and cost.

As a matter of fact, the system characteristics of the data-intensive applications vary tremendously, especially because of varying nature and volume of data leading to no single solution in terms of underlying cluster architecture. The problem becomes more complicated in a public cloud environment or HPC cluster where the application characteristics is not familiar beforehand. In these scenarios when the application characteristics are unfamiliar, as an initial approach it makes sense to invest in a balanced cluster. Hardware vendors as well as cloud vendors (e.g., Amazon, Google, R-HPC, etc.) are investing huge amount of money for this. Millions of dollars are being spent for programs such as, NSF Cloud, XSEDE where several industries and universities collaborate to find such balanced architectures to drive the next generation cloud research.

At this inflection point of HPC landscape, system designers must consider more degrees of freedom for cluster architecture than for existing HPC clusters which focus only on doing calculation at blazing speed. They must address such questions as: how much I/O bandwidth is required per processing core? How much memory is required to optimize the performance and cost? Similar questions can be asked for scaled up and scaled out architecture also, as there is no concrete definition available for these terminologies. These complex performance

A. K. Das, K. Lee, and S.J. Park are with the Department of Computer Science and Electrical Engineering, Center for Computation and Technology Louisiana State University, Baton Rouge, USA

J. Hong and W. Chang are with Samsung Electronics, S. Korea

and economic factors together motivate new design of HPC infrastructure, however, poses several challenges to the system designers who are striving for system balance in terms of processing speed, I/O bandwidth, memory, and the cost of the infrastructure. There is limited analytical studies which address the challenges in developing a cost-effective, balanced distributed cyber infrastructure to analyze large scale scientific big data.

With this motivation, this paper takes an initial attempt to theoretically augment Amdahl's second law for balanced system (in particular, Amdahl's I/O and Memory number) considering both the application characteristics as well as the price of hardware components. It proposes a simple additive model for cost/performance to quantify the system balance between CPU speed, I/O bandwidth, and size of DRAM in terms of software application characteristics and current trend in hardware cost. Our model does not assume any specific software framework or hardware technologies, rather, it provides an easy to use, general guideline for designing a balanced system. However, the outcome of the model (i.e., configuration of an optimal balanced system) suggest to use solid state drives (SSD) instead of hard disk drive (HDD) in a multicore machine to maintain the system balance. Assuming an equal distribution of I/O- and compute work in a data-intensive application, our model suggests that a balanced HPC system needs almost 0.19-GBPS I/O bandwidth, and almost 3-GB of DRAM per GHz of CPU speed using Intel Xeon processor and current price trend of different hardware.

To substantiate our claim, we evaluate three fundamentally different cluster architecture as follows: 1) a traditional HPC cluster, called SuperMikeII (located at LSU, USA) that offers 382 computing nodes each with 16-Xeon cores, 1-HDD, and 32-GB DRAM, 2) a regular datacenter architecture, called SwatIII (located at Samsung, Korea) that has 128 nodes each with 16-Xeon cores, 4-HDD, and 256-GB DRAM and 3) a new MicroBrick-based architecture, called CeresII (also located at Samsung, Korea). Each CeresII node has 6-Xeon cores, 1-NVMe SSD with 2-GBPS I/O bandwidth and 6-GB DRAM per node, thus closely resembles to the optimum produced by our model. We evaluated these three architectures with respect to three different benchmark applications. Two of them are widely used benchmark Hadoop applications (e.g., TeraSort and WordCount) and the other one is our own benchmark genome assembler developed atop Hadoop and Giraph which serves as a good real-world example of a data-, compute- and memory-intensive workload. For all the three benchmark evaluation, CeresII, with the most optimum configuration outperformed the others in terms of both performance and cost/performance when using same amount of processing cores as SuperMikeII and SwatIII across the cluster, and significantly less (almost half) DRAM than SwatIII.

The rest of the paper is organized as follows: Section-II describes the prior work relate to our current effort. Section-?? describes our motivation for this study. In Section-III, we discuss Amdahl's second law in details. Section-IV describes the proposed model. Then, in Section-V shows the details of our experimental testbeds and classifies those using our proposed model. Section-VI describes the evaluation methodology for

these clusters, i.e., the details of the software and benchmark that we used in our evaluation. In Section-VII we discuss the experimental result. Section-?? discusses the limitation in the current model to stimulate discussion and future work. Finally, Section-IX concluded the paper.

## II. RELATED WORK

Numerous studies have been performed evaluating the performance implication of different big data analytic software on different types of hardware infrastructure. We categorize these existing studies into four different classes such as 1) experimental evaluation of different cluster architecture for big data software, 2) simulation of big data software, 3) analytical performance modeling of big data software and 4) system characterization using Amdahl's second law which is the most relevant to our current effort.

### A. Experimental evaluation of different cluster architecture for big data software.

Michael [4] investigated the performance characteristics of the scaled out and scaled up architecture for interactive queries and found better performance using a scaled out cluster. On the other hand, Appuswamy [3] reached an entirely different conclusion in their study. They observed a single scaled up server to perform better than a 8-nodes scaled out cluster for eleven different enterprise-level Hadoop applications including log-processing, sorting, Mahout machine learning, etc. Considering the impact of individual hardware component on overall cluster architecture, several study has also been performed evaluating the impact of different storage media (SSD and HDD) or network on Hadoop. For example, Kang [5] and Wu [6] argued for using SSD to improve the performance of Hadoop cluster with respect to different Hadoop benchmark such as, WordCount, Bayesian, DFSIO, TeraSort etc. Moon [7] showed a significant cost benefit in TeraSort by storing the intermediate Hadoop data in SSD, leaving HDDs to store Hadoop Distributed File System (HDFS [8]) data. A similar result can be found in the study by Li [9] and Krish [10] where SSDs are used to store temporary data to reduce disk contention and HDDs are used to store the HDFS data.

In our previous study [11] we aggregated many of these efforts together and reported the performance result for both in the level of the individual hardware component (network, storage, and memory) as well as their overall all organization. Considering 8 different cluster configurations we provided significant insight on how to deploy SSDs better in Hadoop cluster. We also showed the imbalance between performance and economy that is prevalent in existing scaled up and scaled out architecture. Although we experimentally showed that this imbalance was eliminated in a Samsung MicroBrick-based prototype cluster, in [11] we did not provide any quantitative analysis, either for scaled up/out or for balanced system that can improve performance while staying with the budget. In fact, none of the prior efforts mentioned in this section provide such quantitative analysis which is extremely critical to propose alternative hardware configurations for next generation high performance balanced distributed cyber infrastructure for data- and compute-intensive applications.

### B. Simulation of big data software

Simulation is a widely used for predicting performance of different big data software (mainly Hadoop) and analyzing design trade-offs in different hardware.

At the border level, all simulators generate virtual Hadoop MapReduce jobs and tune different hardware and software parameters in a simulated environment mostly using prior experiences or trial-and-error to optimize the performance of some parts of the Hadoop framework or a given Hadoop job. For example, HSim [12] focuses on the Hadoop job parameters to optimize the performance of the entire job. SimMR [13] on the other hand focus on simulating the Hadoop task scheduling algorithms. MRSim [14] and MRPerf [15] unified all these aspects in a single simulator. They simulate the hardware environment using discrete event simulator packages such as, SimJava, GridSim, NS2, etc. Then they execute the fake Hadoop job on a small subset of data to predict the performance. Given prior experience, simulators can predict hardware alternative, thus capable to save huge cost and effort comparing to experimental evaluation on real hardware discussed before.

However, most of these simulators comes with lots of overhead. They are time and resource consuming and often fail to assess the real system characteristics in the presence of huge volumes of big data. Furthermore, the trial-and-error method of performance optimization considering broad range of available hardware infrastructure with more 200 Hadoop parameters is challenging and most of the time unreliable.

### C. Analytical performance modeling of big data software

Analytical models abstract away several performance parameters and predict the performance of a Hadoop job mainly using single or multi layer queuing network. For example, Viana [16] model the pipeline parallelism of a Hadoop job using queuing network to predict the performance of the job. Wu [17] proposed a layered queue network model to predict the performance of a Hadoop job in a cloud environment. In a recent work Ahn [18] proposed a queuing theory model to predict the performance of Hadoop job atop different storage technologies, i.e, HDD and SSD. Krevat [19] computed the I/O complexity of the Hadoop job for data-intensive applications and proposed a model to quantify the hardware resource wasted by Hadoop.

Unlike simulation, there is no overhead involved in analytical model approach. However, the available alternatives for individual hardware component are huge in number. It is hard to find an optimally balanced architecture in this vast range of alternative hardware components. Thus, both simulation and analytical approaches that model the performance of big data software may work well in selecting alternatives between two or more existing hardware infrastructure (i.e., a small finite search space) however, their capability is limited in proposing new architecture.

### D. System characterization using Amdahl's second law

Till date, the most practical approach to design a balanced system is to follow Amdahl's second law. Computer scientist

Gene Amdahl postulated several design principles for a balanced system, collectively known as Amdahl's Second Law. He stated a balanced system needs 1-bit of sequential I/O per second (Amdahl's I/O number), and 1-byte of memory (Amdahl's memory number) per CPU instruction per second. These design principles for balanced system still holds true even after 50years with some ammendments proposed by Jim Gray in early 2000 [20]. Gray [20] observed that Amdahl's I/O number still holds true but should be measured with respect to the relevant workload. He also observed that Amdahl's memory number is rising from 1 to 4. Gray's ammendments are based on contemporary observation and are not formalized. We will discuss both Amdahl's second law and Gray's ammendment in more details as well as formalize Gray's ammendment later in Section-III.

Amdahl's second law (with Gray's ammendment) can be used for system characterization and proposing balanced system. For example, Bell and Gray [21] classified the existing supercomputers based upon Amdahl's second law to clarify the future roadmap of the HPC architecture. Cohen [22] applied Amdahls second law to the datacenter cluster to study the interplay between processor architecture and network interconnect in a datacenter. Chang [1] used Amdahl's second law to better understand the design implications of data analytic systems by quantifying workload requirements. Szalay [2] using Amdahl's second law proposed a new cluster architecture based on SSD and low power processors (such as, Intel Atom, Zotac etc) to achieve the balance between performance and energy efficiency. In this study, the authors ignored the CPU micro architecture and simplified the Amdahl's I/O and Memory number by dividing the I/O bandwidth (bit/s) and size of DRAM (GB) respectively by CPU-speed (GHz) instead of using instruction per second. The cluster architecture was balanced as the simplified Amdahl's I/O number was close to unity. In a recent study [23] Zheng evaluated a similar architecture as Amdahl's balanced blade for Hadoop applications and showed Atom processor is the system's bottleneck

Our study is different from the above in the following two aspect: Unlike [2], we consider a balanced system as the one which optimizes between both cost and performance. Hence, we did not consider the optimal balance of the system (i.e., the system's I/O and memory ratio to the processor speed) as constants as in Amdahl's I/O and memory number. We express the optimal system balance as a function of both application characteristics and hardware price.

## III. AMDAHL'S SECOND LAW, GRAY'S AMMENDMENT, AND THEIR LIMITATIONS

### A. Original Form of Amdahl's Second Law

Computer scientist Gene Amdahl postulated several design prnciple in late 1960 to make a balanced system. As mentioned earlier, these design principals are collectively known as Amdahls second law which are as follows:

- 1) *Amdahls I/O Law*: A balanced computer system needs one bit of sequential I/O per second per instruction per second. From this point we will mention this law as Amdahls I/O number. Alternatively, Amdahl's I/O

number of a balanced system can be expressed as 0.125 GBPS/GIPS (by changing in conventional units).

- 2) *Amdahls memory Law*: A balanced computer system needs one byte of memory per instruction per second. From this point we will mention this law as Amdahls memory number.

Using the notations in Table-I the original Amdahl's I/O and meory number can be expressed as:

$$\beta_{io}^{opt} = 0.125 \quad (1)$$

$$\beta_{mem}^{opt} = 1 \quad (2)$$

### B. Gray's Ammendment to Amdahl's Second law

Computer scientist Jim Gray reevaluated and ammended Amdahl's second law in the context of modern data engineering. These ammendments are collectively known as Gray's law. The revised laws are as follows:

- 1) *Gray's I/O law*: A system needs 8 MIPS/MBPS I/O, but the instruction rate and IO rate must be measured on the relevant workload. (Sequential workloads tend to have low CPI (clocks per instruction), while random workloads tend to have higher CPI.)
- 2) *Gray's memory law*: The MB/MIPS ratio is rising from 1 to 4. This trend will likely continue.

The underlying implication of the Gray's first ammendment (i.e., Revised Amdahl's I/O Law) is that, it aims for systems with high Amdahl I/O numbers at a given performance level that match the Amdahl I/O numbers of the applications. Regarding Amdahl memory number, in stead of any theoretical justification, Gray put forward a statistics reflecting the contemporary state of cluster architecture.

Using the notations in Table-I Gray's ammendments to original Amdahl's second law can be expressed as follows:

$$\beta_{io}^{opt} = \gamma_{io} \quad (3)$$

$$\beta_{mem}^{opt} = 4 \quad (4)$$

### C. Limitations of Existing Laws

Amdahl's second law for balanced system does not consider the impact of application balance (or, applications' resource requirement) on cluster architecture. Because of the diverse resource requirements, one-size-fit-all design as suggested in the original law (expressed as the constants in the right hand side of Equation-1 and 2) cannot satisfy the different resource balance ratios for a collection of analytic applications.

Gray's law is more realistic in the sense that it consider the impact of application balance (or, applications' resource requirement) on the cluster architecture. However, it is limiting to reflect the interplay between application balance and technology-cost balance. The cost of hardware components has already changed the performance point and will keep on changing as the technology advances.

TABLE I: Notations used in the model and their meaning

$R_{cpu}$	CPU speed of a given system $S$ (GHz)
$R_{io}$	I/O bandwidth of system $S$ (GBPS)
$R_{mem}$	DRAM size of system $S$ (GB)
$W_{cpu}$	Fraction work done by the CPU for a given application $W$
$W_{io}$	Fraction of work done by the disk(s) for $W$
$W_{mem}$	Fraction of work done by DRAM for $W$
$P_{cpu}$	Price per GHz of CPU speed
$P_{io}$	Price per GBPS of I/O bandwidth
$P_{mem}$	Price per GB of DRAM
$\beta_{io}$	System balance between I/O bandwidth and CPU speed for system $S$ ( $= R_{io}/R_{cpu}$ )
$\beta_{mem}$	System balance between DRAM size and CPU speed for system $S$ ( $= R_{mem}/R_{cpu}$ )
$\gamma_{io}$	Application balance between CPU and I/O bandwidth for application $W$ ( $= W_{io}/W_{cpu}$ ). This term quantifies to what extent the application is I/O-intensive or, compute-intensive.
$\gamma_{mem}$	Application balance between CPU and DRAM size for application $W$ ( $= W_{mem}/W_{cpu}$ ). This term quantifies to what extent the application is memory-intensive or, compute-intensive
$\delta_{io}$	Technology-Cost balance of between CPU and I/O bandwidth for system $S$ ( $= P_{io}/P_{cpu}$ )
$\delta_{mem}$	Technology-Cost balance between CPU and DRAM for system $S$ ( $= P_{mem}/P_{cpu}$ )
$\beta_{io}^{opt}$	Optimal system balance between I/O bandwidth and CPU speed (The problem under consideration)
$\beta_{mem}^{opt}$	Optimal system balance between DRAM size and CPU speed (The problem under consideration)

## IV. PROPOSED MODEL FOR SYSTEM BALANCE

### A. Problem Definition

Using the notations described in Table-I the optimal system balance (i.e.,  $\beta_{io}^{opt}$  and  $\beta_{mem}^{opt}$ ) needs to be expressed as a function of application balance (i.e.,  $\gamma_{io}$  and  $\gamma_{mem}$ ) and technology-cost balance (i.e.,  $\delta_{io}$  and  $\delta_{mem}$  in Table-I). Mathematically it can be written as:

$$\beta_{io}^{opt} = f_1(\gamma_{io}, \delta_{io}) \quad (5)$$

$$\beta_{mem}^{opt} = f_2(\gamma_{mem}, \delta_{mem}) \quad (6)$$

### B. Model Assumptions

The model first ignores the same CPU micro architecture. The practical implication of this assumption is that, the optimum system balance (i.e.,  $\beta_{io}$  and  $\beta_{mem}$ ) may vary with different CPU microarchitectures such as, Intel Atom and Xeon, etc. However, the system designers can use this model repeatably for different micro architectures without any loss of generality. On the other hand, the processor micro architecture does not change as frequently as their speed (refer to Intel's *Tick* and *Tock* model for release of new processor technology). Thus, the proposed model significantly reduces the search space for the designers of balanced system.

Second, for simplicity, we assume that the model is additive. That is, we ignore the overlap between work done by CPU, memory and I/O subsystem. This way, the total execution time ( $T_{total}$ ) of an application can be written as:

$$T_{total} = T_{cpu} + T_{io} + T_{mem} \quad (7)$$

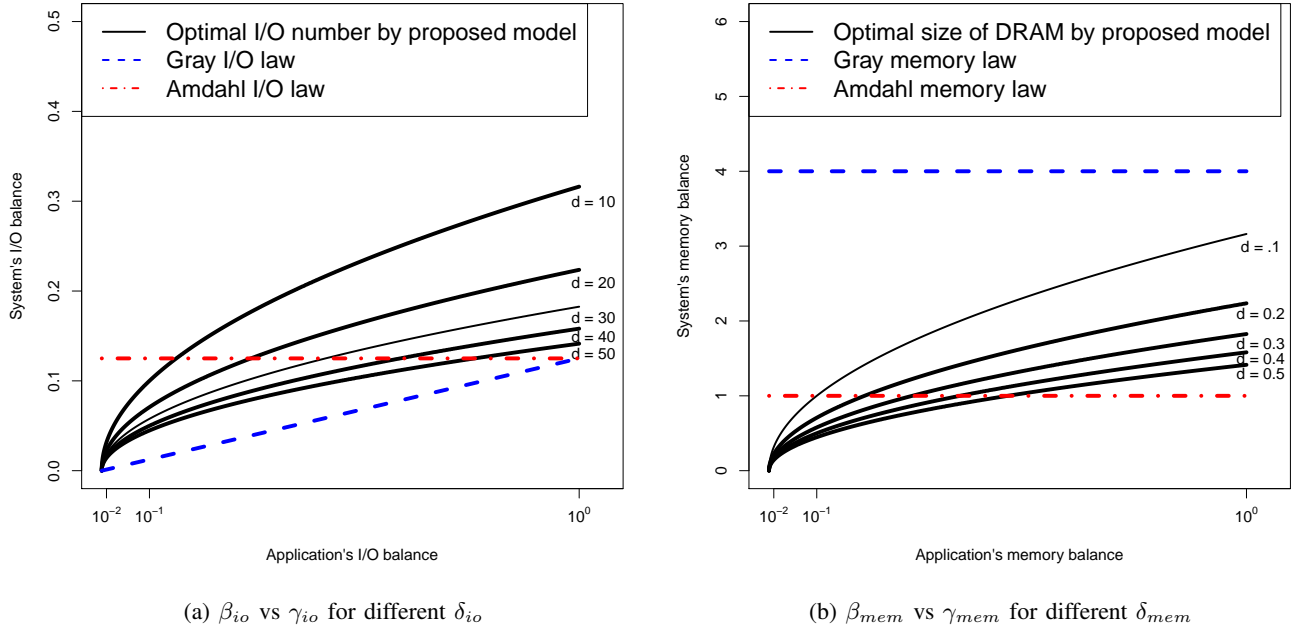


Fig. 1: Change in system's optimum I/O ( $\beta_{io}^{opt}$ ) and memory ( $\beta_{mem}^{opt}$ ) balance for different application balance ( $\gamma_{io}^{opt}$  and  $\gamma_{mem}^{opt}$ ) and technology-cost balance ( $\delta_{io}^{opt}$  and  $\delta_{mem}^{opt}$ )

$$\Rightarrow T_{total} = \frac{W_{cpu}}{R_{cpu}} + \frac{W_{io}}{R_{io}} + \frac{W_{mem}}{R_{mem}} \quad (8)$$

Third, we assume the total cost of the system as the summation of individual cost of CPU, I/O and memory subsystem only. We ignore several constant component such as, base cost, service cost, etc. This way, the total system cost ( $C_{total}$ ) can be written as:

$$C_{total} = C_{cpu} + C_{io} + C_{mem} \quad (9)$$

$$\Rightarrow C_{total} = P_{cpu}R_{cpu} + P_{io}R_{io} + P_{mem}R_{mems} \quad (10)$$

### C. Model Derivation

Assuming the performance as the inverse of the total execution time, the cost/performance (denoted as  $f_{cp}$ ) can be expressed as:

$$f_{cp} = C_{total} \times T_{total} \quad (11)$$

$$\Rightarrow f_{cp} = (C_{cpu} + C_{io} + C_{mem}) \times (T_{cpu} + T_{io} + T_{mem}) \quad (12)$$

$$\begin{aligned} \Rightarrow f_{cp} = & C_{cpu}T_{cpu} + C_{cpu}T_{io} + C_{cpu}T_{mem} \\ & + C_{io}T_{cpu} + C_{io}T_{io} + C_{io}T_{mem} \\ & + C_{mem}T_{cpu} + C_{mem}T_{io} + C_{mem}T_{mem} \end{aligned} \quad (13)$$

Expanding all the time ( $T$ ) and cost ( $C$ ) terms using Equation-8 and 10 respectively, and then, substituting with the notation

used for system balance in Table-I (i.e.,  $\beta_{io}$  and  $\beta_{mem}$ ), Equation-13 can be rewritten as:

$$\begin{aligned} f_{cp} = & P_{cpu}W_{cpu} + \frac{1}{\beta_{io}}P_{cpu}W_{cpu} + \frac{1}{\beta_{mem}}P_{cpu}W_{mem} \\ & + \beta_{io}P_{io}W_{cpu} + P_{io}W_{io} + P_{io}W_{mem} \\ & + \beta_{mem}P_{mem}W_{cpu} + P_{mem}W_{io} + P_{mem}W_{mem} \end{aligned} \quad (14)$$

Assuming a CPU core cannot perform disk and memory operation at the same time, partial differentiation with respect to  $\beta_{io}$  and  $\beta_{mem}$  can separately lead us to the optimum system balance in terms of I/O bandwidth and DRAM size respectively with respect to processing speed.

Partially differentiating with respect to  $\beta_{io}$ ,

$$\frac{\partial f_{cp}}{\partial \beta_{io}} = -\frac{1}{\beta_{io}^2}P_{cpu}W_{io} + P_{io}W_{cpu} \quad (15)$$

For the optimal balance ( $\beta_{io}^{opt}$ ) between CPU speed and I/O bandwidth, Equation-15 should equal to 0. Then, solving for  $\beta_{io}$  and replacing with the workload and cost balance terms mentioned in Table-I we get,

$$\beta_{io}^{opt} = \sqrt{\frac{\gamma_{io}}{\delta_{io}}} \quad (16)$$

Similarly, the optimum balance ( $\beta_{mem}^{opt}$ ) between CPU speed and size of DRAM can be derived as,

$$\beta_{mem}^{opt} = \sqrt{\frac{\gamma_{mem}}{\delta_{mem}}} \quad (17)$$

Equation-16 and 17 show the contribution of application balance and technology-cost balance towards optimal system balance. Gray's law is a special case of our model when the

technology-cost balance equals the inverse of the application balance, i.e., the application's CPU requirement and I/O (or, memory) requirement exactly balance the contemporary cost of the hardware.

Figure-1 shows the change in optimal system balance proposed by our model for different application balance and technology-cost balance. It also compares our model with Amdahl's second law and Gray's laws. It can be noticed that lower is ratio between per GBPS I/O bandwidth cost to per GHz CPU cost higher is the ratio between the system's total I/O bandwidth to system's total CPU speed. That is, for a fixed CPU cost, if the I/O bandwidth cost decreases then, system designers should increase the I/O bandwidth of the system to get the cost optimized performance.

#### D. Notes on Different types of I/O

1) *Sequential vs Random I/O*: It should be noticed that the model did not consider sequential and random I/O bandwidth separately. Designers need to be careful about the application characteristics. An application with many random I/O may find it beneficial to use SSD instead of HDD to align with the optimum I/O balance (i.e.,  $\beta_{io}$ ) produced by the model. For example, Hadoop involves a huge number of random I/O in its shuffle phase. To provide such a huge random I/O bandwidth, SSD can be used.

An underlying assumption of the model is that the aggregated I/O bandwidth of all the disks attached to a compute node is less than or equal to the bandwidth of disk controller. Otherwise, the disk controller can be saturated (an observation by Szalay [2] and our previous work [11]). Because of lower I/O bandwidth of HDD, traditionally disk controllers has never been the source of bottleneck. However, the demand for I/O bandwidth is increasing for big data applications. To meet this demand, it is common to use multiple disks per node. Furthermore, the recent SSDs started producing multiple GBPS of I/O bandwidth. It is designer's responsibility to choose correct hardware so that disk and the disk controller is balanced, i.e., disk controller is neither overprovisioned nor underprovisioned.

2) *Network I/O*: The model did not consider the impact of network interconnect. Hence, we find it important to provide some qualitative information regarding current state of the network architecture in HPC clustre, source of architectural imbalance, and possible solution. Traditional HPC clusters uses a hierarchy of InfiniBand switches with a fat tree model of connectivity. This approach typically uses a layered architecture. The hosts or servers are connected to the bottom layer, called access layer which is then aggregated to an intermediate layer of switches, called edge layer or leaf. The edge layer switches are then connected to the top layer switches, called core layer or spine where the actual bandwidth of the network is scaled. To prevent over subscription, the link speeds got progressively higher from access layer to leaves to spine starting from only few Mbps to several Gbps. This architecture may suffer from bottleneck issues introducing architectural imbalance as the bandwidth of the host adapter is increasing in an outstanding pace (an observation by Cohen

TABLE II: Cost of different hardware components

Hardware components	Cost(\$)	Unit Price
Intel Xeon 64bit 2.6 GHz E5 series (8-cores) processor	1760	\$42/GHz
Intel Xeon D-1541	650	\$54/GHz
Western Digital RE4 HDD (120MBPS), 500GB	132	\$2.25/GBPS/GB
Western Digital VelociRaptor HDD (150MBPS), 600GB	167	\$1.90/GBPS/GB
Samsung 840Pro Series SATAIII SSD (400MBPS), 512GB	450	\$2.25/GBPS/GB
Samsung NVMe SSD PM953 (2GBPS), 950GB	450	\$0.22/GBPS/GB
Samsung DDR3 16GB memory module	159	\$10/GB
32GB 1600MHz RAM (decided by Dell)	140	\$4.37/GB

[22]). Furthermore, to accomodate many servers with fewer switches (thus, reducing the cost of the network) a blocking is used which again limits the performance of big data genomic applications which demand more bandwidth. As an alternative, the simple Clos based architecture is gaining popularity where all the lower layer switches (i.e., leaf) are connected to all the top layer switches (i.e., spine) using a full mesh topology, thus achieving a non-blocking model using inexpensive devices. The data-intensive applications based of Hadoop or Giraph, which transmits huge amounts of data over network in different phases of computation, can use more bandwidth from this all-to-all connection model.

#### E. An Illustrative Example of Applying the Model to build a Balanced Cluster

In this section we demonstrate an example of how to apply our model to build a cost effective balanced cluster. To do that, we chose Intel Xeon as the processor architecture. To reflect today's data-, compute-, and memory-intensive scientific application, we consider the work done by CPU, I/O and memory subsystem (i.e.,  $W_{cpu}$ ,  $W_{io}$ , and  $W_{mem}$ ) is equal for that application. That is, using the notation in Table-I the application balance can be written as:

$$\gamma_{io} = \gamma_{mem} = 1 \quad (18)$$

Table-II shows the current price of different hardware components which are used in our experimental testbeds discussed in Section-V. Column *Unit Price* shows the current price trend for different processor, storage and memory technologies in their corresponding unit. We consider Intel Xeon processor. As it can be seen in Table-II, the cost per MBPS sequential I/O for both HDD and SATA-SSD is almost similar irrespective of change in storage technology provided the same storage space per disk. Whereas, the I/O bandwidth cost started reducing significantly with NVMe SSD. Whereas, the cost per GB of DRAM is increased almost double from DDR2 to DDR3.

We first calculated the average cost of each hardware component from the available list (Table-II) in terms of their corresponding unit price. For example, we have two different Xeon processors: E5-series and D-series as shown in Table-II.

TABLE III: Experimental Testbeds

Cluster name	Processor	CPU Core Speed (GHz)	#Cores /node	CPU Speed (GHz) /node	#Disks /node and type	Seq I/O Band-width /Disk (GBPS)	Total Seq I/O Band-width (GBPS) /node	DRAM /node (GB)	Maximum #Nodes available	$\beta_{io}$	$\beta_{mem}$
SuperMikeII (Traditional Supercomputer)	Two 8-core SandyBridge Xeon	2.6	16	41.6	1-HDD (SATA)	0.15	0.15	32GB	128	0.003	0.77
SwatIII (Regular Datacenter)	Two 8-core SandyBridge Xeon	2.6	16	41.6	4-HDD (SATA)	0.15	0.60	256	16	0.015	6.15
CeresII (MicroBrick-based Hyperscale System)	One 6-core Xeon	2	6	12	1-SSD (NVMe)	2	2	64	40	0.33	5.33

Their unit price is \$42 and \$54. Hence, in this example we selected average unit cost of the processor as \$48. Similarly the cost of DRAM is calculated as \$7.2.

However, calculating the average I/O bandwidth cost is not as simple as in processor or memory. Because, the price of the storage devices is often directed by the capacity (i.e., GB) whereas, the model needs the cost per bandwidth (i.e., GBPS) which again depends on the number of disks. Hence, we need a general methodology to calculate cost per GBPS. To resolve this issue, we first calculated the unit price of a disk in terms of cost per GBPS per GB. Then we selected a common capacity, 1TB for this example that we can use in a server. Then, we multiplied the unit price with capacity (i.e., 1TB) to calculate the cost per GBPS for each of the disk. For example, the cost per GBPS of Western Digital RE4 HDD is \$2304 ( $2.25 \times 1024$ ). Finally, we calculated the average cost of I/O bandwidth. For the disk drives shown in Table-II is the average I/O bandwidth cost is \$1694.

Then, we divided the I/O-cost per GBPS I/O and DRAM-cost per GB by the CPU-cost per GHz to get the technology-cost balance for I/O and memory respectively. Using the notation in Table-I technology-cost balance for this example can be written as:

$$\delta_{io} = 35.30 \quad (19)$$

$$\delta_{mem} = 0.15 \quad (20)$$

Using Equation-16 and replacing the value of  $\gamma_{io}$  and  $\delta_{io}$  from Equation-18 and 19 respectively, we can get the system's I/O balance in terms of GBPS/GHz as:

$$\beta_{io}^{opt} = 0.17 \quad (21)$$

Similarly, using Equation-17, 18, and 20, we can get the system's DRAM balance in terms of GBPS/GHz as:

$$\beta_{mem}^{opt} = 2.7 \quad (22)$$

## V. EXPERIMENTAL TESTBEDS: CHARACTERIZATION USING SYSTEM'S I/O AND MEMORY BALANCE

As mentioned earlier, in our study we evaluate three different cluster architectures. Table-III shows the overview of our experimental testbeds. The first one, SuperMikeII represents a traditional HPC cluster. This LSU HPC cluster oers a total

of 440 computing nodes. However, a maximum 128 can be allocated at a time to a single user. SwatIII represents a regular datacenter. This Samsung datacenter has 128 nodes. However, we used maximum 16 nodes for our experiments. The last one, CeresII is a novel hyperscale system based on Samsung MicroBrick with a maximum of 40 nodes available to us.

### A. Cluster Characterization using proposed model

Figure-1a and 1b shows the change in system's optimum I/O and memory balance respectively with respect to Intel Xeon processor architecture for varying application balance (i.e., application's resource requirement). Based upon this figure characterize each of our experimental testbeds on the basis of  $\beta_{io}$  and  $\beta_{mem}$  and discuss the major pros and cons of each different architecture.

1) *SuperMikeII (Traditional HPC cluster)*: SuperMikeII has two 8-core Intel Sandybridge Xeon processor per node thus oering huge processing power. However, each SuperMikeII node is equipped with only one HDD, thus, limited in terms of I/O bandwidth. Also, each SuperMikeII node has only 32GB DRAM. Thus, SuperMikeII has both,  $\beta_{io} = 0.003$  and  $\beta_{mem} = 0.77$  a magnitude smaller than the optimum produced by our model for a data-, compute- and memory-intensive application as shown in Equation-16 and Equation-17. According to Figure-1 SuperMikeII can provide optimal performance only for compute-intensive applications, precisely for those application with  $\gamma_{io} = x$  and  $\gamma_{mem} = y$ .

2) *SwatIII (Existing Datacenter)*: Unlike SuperMikeII which use only one HDD per node, SwatIII uses 4-HDDs per node using JBOD (Just a Bunch of Disk) conguration while using the same processor (i.e. two 8-core Intel SandyBridge Xeon) as SuperMikeII. Since the I/O throughput increases linearly with number of disks, SwatIII's  $\beta_{io} = 0.015$  is higher than SuperMikeII but lower than the optimum produced by our model for an I/O- and compute-intensive application (Equation-16). On the other hand, each SwatIII node has 256GB DRAM, thus achieve very high value for  $\beta_{mem} = 6.15$ . It is to be noticed that  $\beta_{mem}$  of SwatIII is even higher than the optimum produced by the model (Equation-17). Thus, according to Figure-1 SwatIII can produce optimal performance for moderately I/O-intensive applications and for

memory-intensive applications, precisely when  $\gamma_{io} = x$  and  $\gamma_{mem} = y$ .

3) *CeresII (MicroBrick-based Hyperscale System)*: The last one, CeresII is a novel hyperscale system based on Samsung MicroBricks. It uses one 6-core Intel Xeon processor with a core frequency of 2-GHz. Each MicroBrick (or, simply the computation module) of CeresII consists of a 6-core Intel Xeon processor with a core frequency of 2-GHz, two NVMe-SSD (Samsung) each with an I/O bandwidth of 2GBPS, and 64-GB DRAM.  $\beta_{io}$  of CeresII is 0.32 which is closer to the optimum produced by our model in Equation-16. On the other hand,  $\beta_{mem}$  of each CeresII module is 5.33. Although, it is higher than the optimal, it is significantly less than SwatIII. Thus, CeresII is the most balanced cluster among all the available resources and we expect to get the best cost to performance for today's I/O-, compute- and memory-intensive applications for which  $\gamma_{io} = \gamma_{mem} = 1$ .

## VI. CLUSTER EVALUATION METHODOLOGY

To evaluate the clusters shown in Table-III, we use SuperMikeII as the baseline configuration and compare the performance of other clusters with respect to this. We always use homogeneous configuration across any cluster. One of the nodes in each cluster configuration is always used as Hadoop MasterNode. All other nodes are used as DataNodes.

### A. Software platform: Hadoop and Giraph

To evaluate the relative merits of the clusters with respect to data-intensive scientific applications, first we need a big data analytic platform. The obvious selection is Hadoop which is nowadays the de facto standard for big data analysis and rapidly gaining popularity in scientific computing. We also evaluate the cluster hardware with respect to Giraph which is a large scale graph processing framework developed atop Hadoop. We use Cloudera-Hadoop-2.3.0 and Giraph-1.1.0 for the entire study and use the Cloudera-Manager-5.0.0 for monitoring the system behavior.

### B. Hadoop Configurations Overview

Since our goal is to evaluate the underlying hardware components and their, we avoid any unnecessary change in the source code of Hadoop or Giraph. To evaluate the balance between the hardware components, we set up the Hadoop parameters such that the application can make use of the maximum available processing speed, I/O bandwidth and DRAM. However, it is worthy to mention here, due to limited I/O bandwidth in SuperMikeII (traditional HPC cluster with only 1-HDD per node) we could launch only 8 concurrent YARN containers (i.e., 8 concurrent map/reduce tasks) i.e., half of the total number of cores to get the optimal performance. Appendix-A shows a brief description of our Hadoop configurations.

### C. Benchmark Application Characteristics

Table-IV summarizes the details of the benchmark applications, the data size handled by each of these applications and

their corresponding characteristics. We used three benchmark application for our evaluation: TeraSort, WordCount, and a real world genome assembly application. A brief description of each of these applications are as follows:

1) *TeraSort*: The first one is TeraSort, a standard benchmark for any data processing framework. The Hadoop version of TeraSort samples and partitions the input data in its map phase based upon only first few characters. The reduce phase then uses quicksort algorithm to sort each of the partition locally. The map phase of TeraSort is CPU-intensive as it reads only the first few characters of each row in the input dataset to generate the key (called sample-key) for each data. However, based upon the data size reduce phase can be severely I/O-intensive. The entire dataset needs to be transferred to different reducers and each reducer needs to read each row at least once to sort the data set. Again, for huge volume of data partition, that cannot fit in memory, an external sort may be applied making the reduce phase severely I/O intensive.

2) *WordCount*: The second one is WordCount, another widely used Hadoop benchmark. The map phase of WordCount parses the input data set line by line to extract each word. Each word is emitted with a 1 as its initial count which is then summed up in the reduce phase to output its frequency. Since both the map and reduce phase read the entire data set sequentially just once, both the map and reduce phase of WordCount is CPU-intensive.

3) *Genome Assembly*: The third one is a genome assembly application that we developed based on Hadoop and Giraph to address the challenges in large scale genome assembly that recently made its way to the forefront of big data challenges. The first phase uses Hadoop to scan the genome data set to filter out the lines containing only nucleotide characters (i.e., A, T, G, and C) and then divides each of those lines into smaller fragments of length  $k$  known as  $k$ -mer. Map phase emits each two subsequent  $k$ -mers as key-value pair representing a vertex and an edge from it. The reduce phase then aggregates all the edges emitted from each vertex to write the entire graph structure to HDFS. The second phase of the assembly application uses Giraph to compress each of the linear chains in the graph structure to one vertex and remove the tip and bubble structure in the graph (included because of sequencing error). After removing the tips and bubble structure the graph may have some new linear chains which are compressed again. Similarly, new tip and bubble structures are removed again. So, the process continues incrementally as a series of Giraph jobs until there is no tip or bubble in the graph. Appendix-B shows the overview of the algorithms used in our genome assembly application.

### D. Data Size

For both TeraSort and WordCount benchmark, we ran TeraGen program and generated another 1TB random dataset as the input for these benchmarks. Both TeraSort and WordCount generate almost similar amount data in the intermediate shuffle phase as their input. The output data size of TeraSort is also same as its input (i.e., 1TB in this work) whereas, the output of WordCount may vary based upon the frequency of different



TABLE IV: Data size for different benchmark applications

Job name	Job Type	Input	Final output	# jobs	Shuffled data	HDFS Data	Application Characteristics
Terasort	Hadoop	1TB	1TB	1	1TB	1TB	Map: CPU-intensive, Reduce: I/O-intensive
Wordcount	Hadoop	1TB	1TB	1	1TB	950GB	Map and Reduce: CPU-Intensive
Graph Construction (large size human genome)	Hadoop	452GB (2-billion reads)	3TB	2	9.9TB	3.2TB	Map and Reduce: CPU- and I/O-intensive
Graph Simplification (large size human genome)	Series of Graph jobs	3.2TB (1483246722 vertices)	3.8GB (3032297 vertices)	15	-	4.1TB	Memory-Intensive

TABLE V: Resources available for each cluster architecture when scaling with respect to cost. The cost of 15 datanodes of SuperMikeII is used as baseline. This is used for TeraSort and WordCount benchmark

Cluster Configurations	SuperMikeII	SwatIII	CeresII
Total cost (\$)	57060	57060	57060
Cost/DN (\$)	3804	6911	2700
#Nodes	16	9	23
Total processing speed (GHz)	624.00	330.83	300.00
Total I/O bandwidth (GBPS)	2.25	4.80	50.00
Total storage space (TB)	7.50	16.00	25.00
Total DRAM Size (TB)	0.48	2.00	1.60

words in the randomly generated dataset. In our case it was almost 950GB.

For the genome assembly benchmark application, in this work we use a large size human genome dataset (452GB). The corresponding graph size is 3.2TB. The Hadoop stage of the assembly application is severely shuffle intensive. The temporary shuffle data is almost 21-times more than the input size. The Human genome is openly available in NCBI website with accession number SRX016231<sup>1</sup>.

## VII. RESULT AND DISCUSSION

### A. Evaluating the Cost/Performance of Different Cluster with TeraSort and WordCount

To show the balance between performance and economy we keep the total cost same across all the clusters. Table-II shows the available resources for all the three cluster architectures available to us i.e. SuperMikeII, SwatIII and CeresII while keeping the total cost same across all the clusters. We used the cost of 15 datanodes of SuperMikeII as the baseline. Then, we divided this baseline-cost by the cost of each node in SwatIII and CeresII to count the number of nodes to be used in these two different clusters.

We ran the first two applications described in Table-IV in all the three cluster configurations and measured their execution time. All results are the means of at least 3 runs of each job on each configuration. Figure-VI shows the results normalized

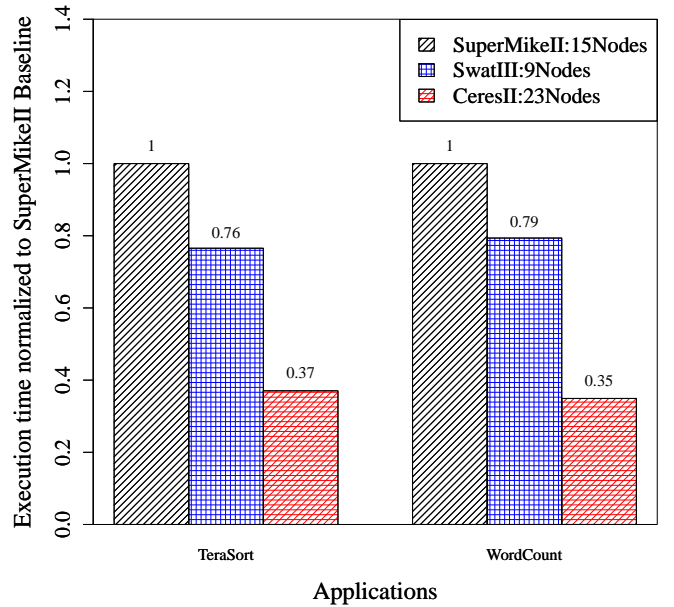


Fig. 2: Execution time of TeraSort and WordCount over different cluster architecture keeping the total cost of each cluster same

to the SuperMikeII baseline. That is, the execution time of SuperMikeII is always 1. We see that CeresII being closer to the optimal produced by our model performs surprisingly well:

- 1) Comparing to SuperMikeII baseline, for both TeraSort and WordCount benchmark application, CeresII shows almost 65% improvement.
- 2) Comparing to SwatIII, CeresII shows almost 50% improvement in execution time for TeraSort and WordCount.

1) *System Characteristics:* Figure-3 and ?? shows the total number of CPU instructions per second (IPS) and I/O bandwidth across different cluster architecture while keeping the total cost of the cluster same. For both TeraSort and WordCount, the peak IPS of SuperMikeII is significantly higher than both SwatIII and CeresII. However, the average I/O bandwidth of SuperMikeII is significantly less than the

<sup>1</sup>[http://www.ncbi.nlm.nih.gov/sra/SRX016231\[accn\]](http://www.ncbi.nlm.nih.gov/sra/SRX016231[accn])

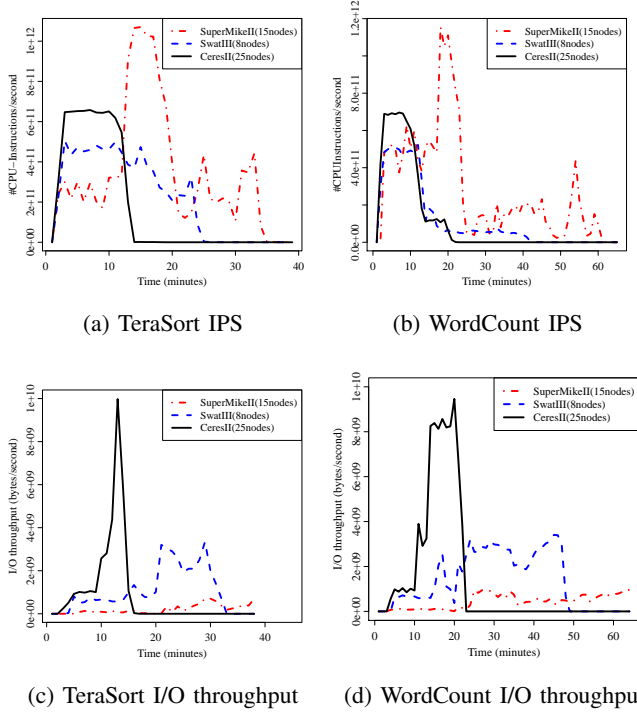


Fig. 3: System Characteristics of TeraSort and WordCount benchmark on different cluster architecture

other two leading to significantly high I/O wait which resulting in lower average IPS. This huge difference between peak and average IPS is because of the extremely low value of  $\beta_{io}$  in SuperMikeII.

SwatIII cluster shows a better balance between CPU and I/O than SuperMikeII. However, because of its high size of DRAM, it gives up processing speed with respect to CeresII (128 core in SwatIII vs 150 cores in CeresII) when the total cost of the cluster is same.

CeresII cluster is the optimally balanced in terms of CPU speed and I/O bandwidth with a  $\beta_{io}$  value of 0.19. Figure-3a and 3b shows that CeresII achieves higher Peak IPS than SwatIII but less than SuperMikeII for both TeraSort and WordCount application. On the other hand, CeresII achieves significantly higher I/O bandwidth than both the cluster keeping the peak and average IPS of the application almost similar and finally leading to the lowest execution time for any of the applications.

### B. Evaluating Different Cluster for Large Size Human Genome Assembly

to assemble the large human genome (452GB) we used the maximum available resources in each of the clusters to accommodate the huge amount of shuffled data (9.9TB) as well as the graph data (3.2TB). That is, we used all 128 nodes of SuperMikeII, 16 nodes of SwatIII, and 40 nodes of CeresII for this application. Figure-4a and 4b shows the execution time and the cost/performance respectively for the Hadoop and Giraph stage of the human genome assembly pipeline. As shown in the model, we consider the performance as the

TABLE VI: Maximum available resources in each cluster architecture (used for large human genome assembly)

Cluster Configurations	SuperMikeII	SwatIII	CeresII
Total cost (\$)	486912	110576	10800
Cost/DN (\$)	3804	6911	2700
#Nodes	128	16	40
Total processing speed (GHz)	5324.8	665.6	480.00
Total I/O bandwidth (GBPS)	19.2	9.60	80.00
Total storage space (TB)	7.50	16.00	25.00
Total DRAM Size (TB)	4	4	2.5

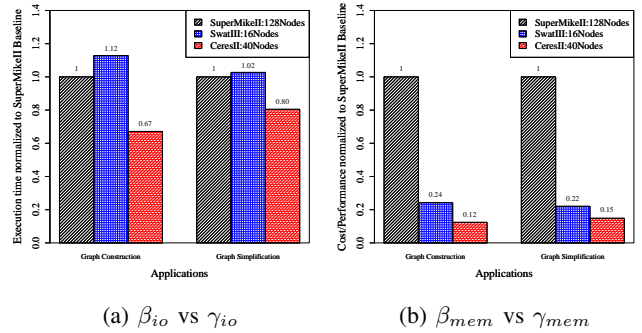


Fig. 4: Performance of different cluster for large size human genome assembly

inverse of the execution time and multiplied the execution time with the total cost of each cluster to get the corresponding cost/performance. All data is again normalized to SuperMikeII baseline. The results are as follows:

- 1) CeresII even with only 9% processing speed than SuperMikeII across the cluster, outperforms SuperMikeII by more than 30% in the Hadoop stage of the assembly. In the Giraph stage, the performance gain is almost 20%. In terms of cost/performance, the corresponding gains are 88% and 85% respectively.
- 2) Comparing to SwatIII, the processing power of CeresII is 72%. However, due to the optimal architectural balance, CeresII outperforms SwatIII by almost 50% in the Hadoop stage. The corresponding improvement in Giraph stage is 20%. In terms of cost/performance, CeresII shows almost 50% and 30% improvement over SwatIII for Hadoop and Giraph stage respectively.

### VIII. LIMITATIONS OF THE MODEL

As an initial approach our main focus was on simplicity. To keep the model simple, we gave up several details and micro architecture in CPU, storage and memory subsystems. It is always possible to include more subtle parameters such as CPU multi-threading, number of I/O operations per second and I/O latency, Network interconnect across cluster, memory-bandwidth and latency, etc. to the existing additive model to come up with better and more accurate equation for system balance. We also ignored the overlap between work done

by CPU and disk which severely depends on the software framework and application's I/O characteristics and its nature of the parallelism (partitioned and pipelined). Building a better, accurate model is the part of our future work.

## IX. CONCLUSION AND FUTURE WORK

Big data needs big resources. As the nature of scientific computing is changing from compute-centric to data-centric, it is obvious that providing more resources (more processing speed, I/O bandwidth, DRAM, etc.) will provide more performance. So, the major challenge is now in providing expected performance in reduced cost. To the best of our knowledge, this paper makes an initial attempt to analytically resolve the performance and cost conundrum prevalent in big data research.

We propose a simple additive model to express the optimal system balance in terms of application characteristics and hardware price to minimize the cost/performance. We also verified our claim experimentally that the cluster that resembles to the optimal produced by our model, i.e., CeresII outperforms an existing HPC cluster and a regular data center architecture by several magnitude for different types and size of applications.

The proposed model also has implications for the way current big data HPC clusters or data center clusters are provisioned. Instead of evaluating the clusters with respect to scaled up and scaled out architecture we argued for a system balance. We believe that if proper balance is maintained between different hardware components, both scaled out and scaled up cluster can provide good performance while maintaining the budget. However, because of resource limitation we could not validate our claim and thus, obviously one of the future directions of our research.

### APPENDIX A

#### HADOOP AND GIRAPH: PROGRAMMING MODEL AND CONFIGURATION

##### A. Programming model of Hadoop and Giraph

Hadoop and Giraph were originated as the open-source counterpart of Google's MapReduce [24] and Pregel [25] respectively. Both the software read the input data from the underlying Hadoop Distributed File System (HDFS) in the form of disjoint sets or partitions of records. Then, in the MapReduce abstraction, a user-defined map function is applied to each disjoint set concurrently to extract information from each record in the form of intermediate key-value pairs. These key-value pairs are then grouped by the unique keys and shuffled to the reducers. Finally, a user-defined reduce function is applied to the value-set of each key, and the final output is written to the HDFS. The MapReduce framework enables data- and compute-intensive applications to run large volume of distributed data sets over distributed compute nodes with local storage. On the other hand, Giraph uses the Bulk Synchronous Parallel model [26] where computation proceeds in supersteps. In the first phase of a superstep, Giraph leverages Hadoop-mappers when a user-defined vertex-program is applied to all the vertices concurrently. In the end of each superstep, each vertex can send a message to other vertices

to initiate the next superstep. Alternatively, each vertex can vote to halt. The computation stops when all the vertices vote to halt unanimously in the same superstep. Giraph enables memory- and compute-intensive applications to upload data into distributed memories over different compute nodes.

##### B. Hadoop configuration and optimization

**Number of concurrent YARN containers:** After performing rigorous testing, we observed, that for SuperMikeII and SwatIII-Basic-HDD (1-HDD/DN cases), 8-containers/DN (i.e., half of total cores/node) show the best result. For any other cluster, number of concurrent containers per datanode is kept equal to the number of cores per node.

**Amount of memory per container and Java-heap-space:** In each node in any node of any cluster, we kept 10% of the memory for the system use. The rest of the memory is equally divided among the launched containers. The Java heap space per worker is always set to lower than memory per container as per normal recommendation.

**Total number of Reducers:** Based on the observation of job profiles, we observed that 2-times of reducers than number of concurrent containers produce good performance in general.

**Giraph workers:** The number of Giraph workers is set according to the number of concurrent YARN containers.

**Other Giraph parameters:** We always used enough memory to accommodate the graph structure in memory and always avoided using the out-of-core execution feature of Giraph, which writes huge data to the disk. We also avoided using the checkpoint feature for the same reason.

### APPENDIX B

#### GENOME ASSEMBLY ALGORITHMS

The motivation behind selecting genome assembly application is that, the high throughput next generation DNA sequencing machines (e.g., Illumina Genome Analyzer) has outpaced Moore's law and started producing a huge amount of short read sequences typically in the scale of several Gigabytes to Terabytes. Furthermore, the size of the de Bruijn graph built from these vast amount of short reads may be a magnitude higher than the reads itself making the entire assembly pipeline severely data-intensive.

De novo genome assembly refers to the construction of an entire genome sequence from a huge amount of small, overlapping and erroneous fragments called short read sequences while no reference genome is available. The problem can be mapped as a simplified de Bruijn graph traversal [27]. We classified the de novo assembly in two stages as follows: a) Hadoop-based de Bruijn graph-construction and b) Giraph-based graph-simplification. In this section, we provide a brief overview of each stage of the assembler.

##### A. Hadoop-based De Bruijn graph-construction (I/O- and compute-intensive workload)

After filtering the actual short reads (i.e., the line containing only nucleotide characters *A*, *T*, *G*, and *C*) from a standard fastq file, an extremely shuffle-intensive Hadoop job creates

the graph from these reads. the Hadoop map task divides each read into several short fragments of length  $k$  known as  $k$ -mers. Two subsequent  $k$ -mers are emitted as an intermediate key-value pair that represents a vertex and an edge (emitted from that vertex) in the de Bruijn graph. The reduce function aggregates the edges (i.e the value-list) of each vertex (i.e., the  $k$ -mer emitted as key) and, finally, writes the graph structure in the HDFS in the adjacency-list format. Based upon the value of  $k$  (determined by biological characteristics of the species), the job produces huge amount of shuffled data. For example, for a read-length of 100 and  $k$  of 31 the shuffled data size is found to be 20-times than the original fastq input. On the other hand, based upon the number of unique  $k$ -mers, the final output (i.e., the graph) can vary from 1 to 10 times of the size of the input.

### B. Giraph-based Graph Simplification (memory- and compute-intensive workload)

The large scale graph data structure produced by the last MapReduce stage is analyzed here. This stage consists of a series of memory-intensive Giraph jobs. Each Giraph job consists of three different types of computation: compress linear chains of vertices followed by removing the tip-structure and then the bubble-structure (introduced due to sequencing errors) in the graph. Giraph can maintain a counter on the number of supersteps and the master-vertex class invokes each type of computation based on that.

### ACKNOWLEDGMENT

The authors would like to thank...

### REFERENCES

- [1] J. Chang, K. T. Lim, J. Byrne, L. Ramirez, and P. Ranganathan, "Workload diversity and dynamics in big data analytics: implications to system designers," in *Proceedings of the 2nd Workshop on Architectures and Systems for Big Data*. ACM, 2012, pp. 21–26.
- [2] A. S. Szalay, G. C. Bell, H. H. Huang, A. Terzis, and A. White, "Low-power amdahl-balanced blades for data intensive computing," *ACM SIGOPS Operating Systems Review*, vol. 44, no. 1, pp. 71–75, 2010.
- [3] R. Appuswamy, C. Gkantsidis, D. Narayanan, O. Hodson, and A. Rowstron, "Scale-up vs scale-out for hadoop: Time to rethink?" in *Proceedings of the 4th annual Symposium on Cloud Computing*. ACM, 2013, p. 20.
- [4] M. Michael, J. E. Moreira, D. Shiloach, and R. W. Wisniewski, "Scale-up x scale-out: A case study using nutch/lucene," in *Parallel and Distributed Processing Symposium, 2007. IPDPS 2007. IEEE International*. IEEE, 2007, pp. 1–8.
- [5] Y. Kang, Y.-s. Kee, E. L. Miller, and C. Park, "Enabling cost-effective data processing with smart ssd," in *Mass Storage Systems and Technologies (MSST), 2013 IEEE 29th Symposium on*. IEEE, 2013, pp. 1–12.
- [6] D. Wu, W. Luo, W. Xie, X. Ji, J. He, and D. Wu, "Understanding the impacts of solid-state storage on the hadoop performance," in *Advanced Cloud and Big Data (CBD), 2013 International Conference on*. IEEE, 2013, pp. 125–130.
- [7] S. Moon, J. Lee, and Y. S. Kee, "Introducing ssds to the hadoop mapreduce framework," in *Cloud Computing (CLOUD), 2014 IEEE 7th International Conference on*. IEEE, 2014, pp. 272–279.
- [8] D. Borthakur, "The hadoop distributed file system: Architecture and design," *Hadoop Project Website*, vol. 11, no. 2007, p. 21, 2007.
- [9] B. Li, E. Mazur, Y. Diao, A. McGregor, and P. Shenoy, "A platform for scalable one-pass analytics using mapreduce," in *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*. ACM, 2011, pp. 985–996.
- [10] K. Krish, A. Khasymski, G. Wang, A. R. Butt, and G. Makkar, "On the use of shared storage in shared-nothing environments," in *Big Data, 2013 IEEE International Conference on*. IEEE, 2013, pp. 313–318.
- [11] A. K. Das, S.-J. Park, J. Hong, and W. Chang, "Evaluating different distributed-cyber-infrastructure for data and compute intensive scientific application," in *Big Data (Big Data), 2015 IEEE International Conference on*. IEEE, 2015, pp. 134–143.
- [12] Y. Liu, M. Li, N. K. Alham, and S. Hammoud, "Hsim: a mapreduce simulator in enabling cloud computing," *Future Generation Computer Systems*, vol. 29, no. 1, pp. 300–308, 2013.
- [13] A. Verma, L. Cherkasova, and R. H. Campbell, "Play it again, simmr!" in *Cluster Computing (CLUSTER), 2011 IEEE International Conference on*. IEEE, 2011, pp. 253–261.
- [14] S. Hammoud, M. Li, Y. Liu, N. K. Alham, and Z. Liu, "Mrsim: A discrete event based mapreduce simulator," in *Fuzzy Systems and Knowledge Discovery (FSKD), 2010 Seventh International Conference on*, vol. 6. IEEE, 2010, pp. 2993–2997.
- [15] G. Wang, A. R. Butt, P. Pandey, and K. Gupta, "A simulation approach to evaluating design decisions in mapreduce setups," in *MASCOTS*, vol. 9. Citeseer, 2009, pp. 1–11.
- [16] E. Vianna, G. Comarela, T. Pontes, J. Almeida, V. Almeida, K. Wilkinson, H. Kuno, and U. Dayal, "Analytical performance models for mapreduce workloads," *International Journal of Parallel Programming*, vol. 41, no. 4, pp. 495–525, 2013.
- [17] X. Wu, Y. Liu, and I. Gorton, "Exploring performance models of hadoop applications on cloud architecture," in *Proceedings of the 11th International ACM SIGSOFT Conference on Quality of Software Architectures*. ACM, 2015, pp. 93–101.
- [18] S. Ahn and S. Park, "An analytical approach to evaluation of ssd effects under mapreduce workloads," *JOURNAL OF SEMICONDUCTOR TECHNOLOGY AND SCIENCE*, vol. 15, no. 5, pp. 511–518, 2015.
- [19] E. Krevat, T. Shiran, E. Anderson, J. Tucek, J. J. Wylie, and G. R. Ganger, "Applying performance models to understand data-intensive computing efficiency," DTIC Document, Tech. Rep., 2010.
- [20] J. Gray and P. Shenoy, "Rules of thumb in data engineering," in *Data Engineering, 2000. Proceedings. 16th International Conference on*. IEEE, 2000, pp. 3–10.
- [21] G. Bell, J. Gray, and A. Szalay, "Petascale computations systems: Balanced cyberinfrastructure in a data-centric world," 2005.
- [22] D. Cohen, F. Petrini, M. D. Day, M. Ben-Yehuda, S. W. Hunter, and U. Cummings, "Applying amdahl's other law to the data center," *IBM Journal of Research and Development*, vol. 53, no. 5, pp. 5–1, 2009.
- [23] D. Zheng, A. Szalay, and A. Terzis, "Hadoop in low-power processors," *arXiv preprint arXiv:1408.2284*, 2014.
- [24] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [25] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel: a system for large-scale graph processing," in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. ACM, 2010, pp. 135–146.
- [26] T. Cheatham, A. Fahmy, D. Stefanescu, and L. Valiant, "Bulk synchronous parallel computing paradigm for transportable software," in *Tools and Environments for Parallel and Distributed Systems*. Springer, 1996, pp. 61–76.
- [27] P. A. Pevzner, H. Tang, and M. S. Waterman, "An eulerian path approach to dna fragment assembly," *Proceedings of the National Academy of Sciences*, vol. 98, no. 17, pp. 9748–9753, 2001.