

Scientific Data Infrastructure

Chetan Anand

Computer Science Engineering
Amc Engineering College,
Bannerghatta, Bangalore, India.

Abstract:- Data in computer world relates to information of something or someone. This data has to be reliable, big in volume, efficient but data in computer world is in 0's and 1's so to understand these data we need a special method. Big data is one such computational method which is designed to overcome this problem. The paper focuses on description of big data and the difficulties faced by the various researcher communities to convert the data, share, implementation the data using a digital infrastructure and future scope of this scientific digital data infrastructure. The paper introduces SDI generic architecture model and how SDLM and SDI can be naturally implemented using cloud.

Keywords:- Big data, scientific data infrastructure (SDI), scientific data life cycle management (SDLM), cloud.

I. INTRODUCTION

Big data is a term that describes the large volume of data (both structured and unstructured) that inundates a business on a day-to-day basis. But it's not the amount of data that's important. It's what organizations do with the data that matters. Big data can be analyzed for insights that lead to better decisions and strategic business moves at global level. Big data is being used from 2000 by first Seisint Inc. then in 2004 LexisNexis acquired Seisint Inc. in 2011 Apache V2.0 came in market till now only publically available platform which integrated HPC AND QUANTCAST file system. In 2001 Doug Laney characterized big data on the basis of few parameters those are:-

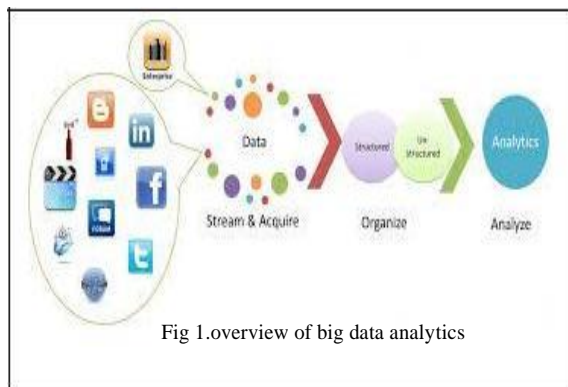


Fig 1. overview of big data analytics

- **Volume:** Organizations collect data from a variety of sources including business transactions, social media and information from sensor or machine-to-machine data. In the past, storing it would've been a problem but new technologies (such as Hadoop) have eased the burden.
- **Velocity:** Data streams in at an unprecedented speed and must be dealt with in a timely manner.
- **Variety:** Data comes in all types of formats from structured (numeric data to traditional databases) to unstructured (text documents, email, video, audio, stock ticker data and financial transactions).
We consider two additional dimensions when it comes to big data: -
- **Variability:** In addition to the increasing velocities and varieties of data, data flows can be highly inconsistent with periodic peaks. Daily, seasonal and event-triggered peak data loads can be challenging to manage.
- **Complexity:** Today's data comes from multiple sources, which makes it difficult to link, match and transform data across systems. However, it's necessary to connect and correlate relationships, hierarchies and

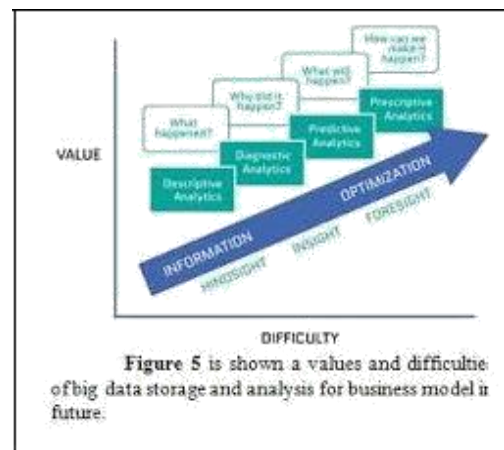


Figure 5 is shown a values and difficultie of big data storage and analysis for business model in future.

Multiple data linkages or your data can quickly spiral out of control [6].

II. PROBLEMS IDENTIFIED IN BIG DATA

Researchers have found out some of the difficulties on working with big data. These difficulties are:
 1. Privacy and security
 2. Data access and sharing of information
 3. Analytical challenge
 4. Human resources and manpower
 5. Distribution of data
 these limitations are overcome by the scientific data infrastructure (SDI) model.

III. SCIENTIFIC DATA INFRASTRUCTURE (SDI) AND (SDLM)

Once the data is published, it is essential to allow the other scientist to be able to validate and reproduce the data that they are interested in and possibly contribute with new results. Another aspect to take in consideration is that the data is distributed both on processing and collection side linking these data is difficult. However as we anticipate the scale of published data semantic becomes necessary condition for efficient reuse of published data. The European commission's initiative to support open access to scientific data from publicly funded projects suggests

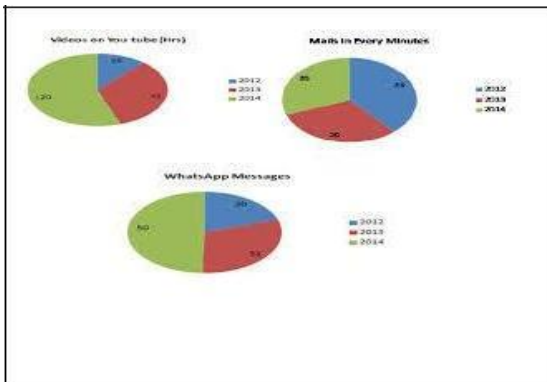


Fig 2. Amount of huge data generated

- PID-persistent data ID
- ORCID-open Researcher and contributor identifier.

Therefore, the researchers introduced scientific Data Lifecycle Management (SDLM) MODEL based on the result of analyses in different scientific communities. The generic scientific data lifecycle includes a number of stages: research project, data collection, data processing, publishing research results, discussion, feedback, archiving. New SDLM requires data storage and preservation at all stages should allow data re-use/re-purposing and secondary research on the processed data and published results. However this is only possible if the full data identification, cross-reference and linkage are implemented in SDI. The SDLM should support data security and access control to scientific data during their

lifecycle: data acquisition, initial data filters, specialist processing; research data storage and secondary data mining, data and research information archiving.

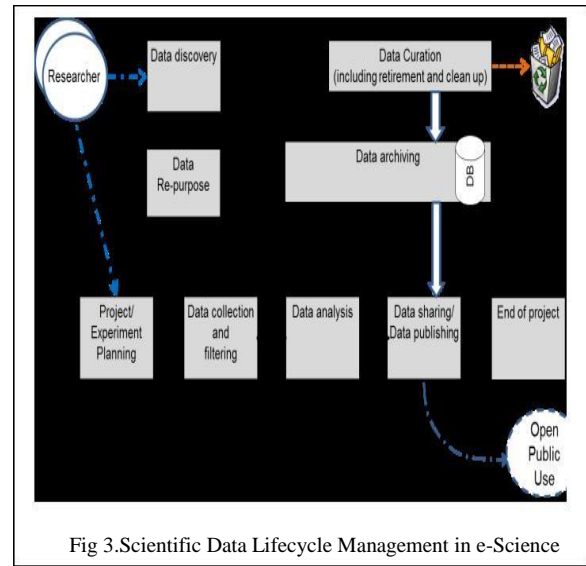


Fig 3. Scientific Data Lifecycle Management in e-Science

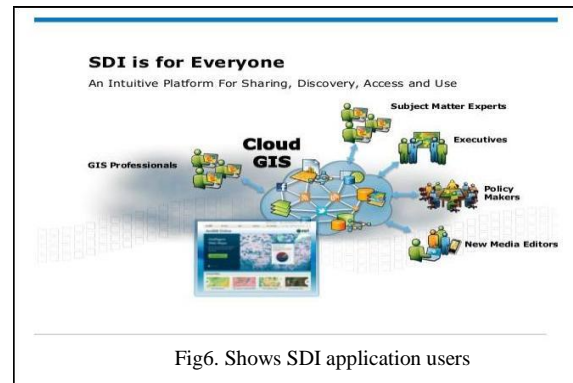


Fig6. Shows SDI application users

We also propose SDI Architecture for e-Science (e-SDI) is illustrated in Figure 4. It contains the following layers:

Layer D1:- Network infrastructure layer represented by the general purpose Internet infrastructure and dedicated network infrastructure.

Layer D2:- Datacenters and computing resources/facilities.

Layer D3:- Infrastructure virtualization layer that is represented by the Cloud/Grid infrastructure services and middleware supporting scientific platforms deployment and operation.

Layer D4:- (Shared) Scientific platforms and instruments specific for different research areas.

Layer D5:- Federation and Policy layer that includes federation infrastructure components, including policy and collaborative user groups support functionality.

Layer D6:- Scientific applications and user portals/clients.

The three cross-layer planes are defined: Operational Support and Management System; Security plane; and Metadata and Lifecycle Management.

The dynamic character of SDI and its support of distributed multi-faceted communities are supported by the dedicated layers: D3 – Infrastructure Virtualisation layer that typically uses modern cloud technologies; and D5 – Federation and policy layer that incorporates related federated infrastructure management and access technologies.

enterprise model can be mapped to cloud based services and after that deployed and operated at any instant inter-cloud infrastructure.it contains cloud infrastructure segments IaaS (VR3-VR5) and PaaS (VR6-VR7), separate virtualized resources (VR1-VR2), two interacting campuses A and B and interconnecting them network infrastructure that in many cases may need to use dedicated network links for guaranteed performance [2][3].

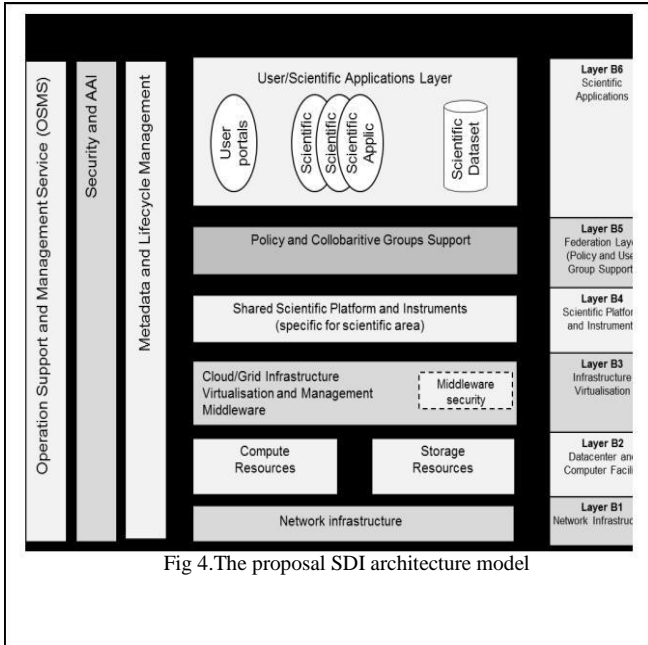


Fig 4.The proposal SDI architecture model

V. ADVANTAGES AND DISADVANTAGES OF MODEL

- A. **DISADVANTAGE:** - For interaction between interacting campuses there is need for dedicated link for which overall infrastructure management and individualservices is required which is typically out of scope of this model [8].
- B. **ADVANTAGE:** -
 1. Solves the problem of distribution of data.
 2. Provides security and access control touser.
 3. Enables sharing of data and re-usability of data.
 4. Provides future scope of improvement on data and model.
 5. Easy to work for user [1].

VI. LITERATURE SURVEY

Author name	Title paper	Highlights of paper
Yuri demchenko, Zhiming Zhao, Paola Grosso,Adianto Wibisono,Cees de Laat	Addressing big data challenges for scientific data infrastructure	SDLM and SDI infrastructure model
Yuri demchenko, Canh Ngo, Peter Membrey	Architecture framework and components for big data ecosystem	Data Models and Structures , Data Management and Big Data Lifecycle,Cloud Based Infrastructure Services for BDI

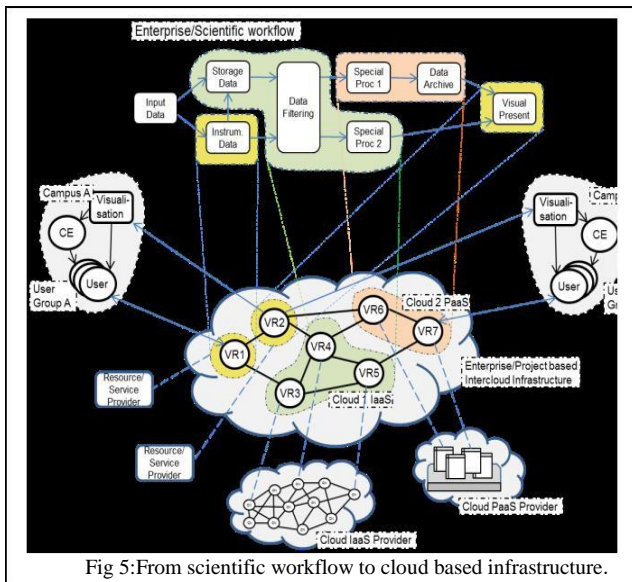


Fig 5:From scientific workflow to cloud based infrastructure.

Now let us see the cloud implementation of the SDI architecture model:-

The main goal of the enterprise or scientific workflow and operational procedures related to processes monitoring and data processing. Cloud technologies simplify the building of such infrastructures and provision it on demand.The figure 5 illustrates how an scientific or

VII. FUTURE RESEARCH ANDDEVELOPMENT

The future research and development will include further SDLM definition, e-SDI and ACAI components definition and development with focus on infrastructure components of e-SDI. Special attention will be given to defining the whole cycle of the provisioning SDI services on-demand specifically tailored to support instant scientific workflows using cloud IaaS and PaaS platforms. This

research will be also supported by development of the corresponding cloud and Inter Cloud architecture framework to support Big Data e-Science processes and infrastructure operation.

VIII. CONCLUSION

Here we can conclude the following things from the paper Big data computes data in huge volume, variety and complex data whose relation can be with anything and gives profit to the organization.

- The limitations of big data is privacy , distribution, analytics which make the model less effective in order to overcome that we have SDLM model with SDI architecture.
- The SDI architecture is the basis of the implementation of any big data model with data security, access, distribution of the data, controls the flow of data and storage of data.
- The SDI model of cloud implementation is best for distributed data and makes linking of data possible.
- The SDLM model takes the overview of the architecture of storage of data and proposes the SDI architecture.
- Although this model has flow and future research is still on progress the SDI model can still fulfill the current needs of big data computation.
- This architecture is currently used by Intel and Google and other companies.

IX. REFERENCES

- [1] <https://tdwi.org/Articles/2015/10/27/5-Best-Practices-Implementing-Big-Data.aspx>
- [2] Study on AAA platforms for scientific data (<https://confluence.terena.org/download/attachements/30474266/AAA-study-report-0907.pdf>)
- [3] EGI federated cloud task force. <Http://www.egi.eu/infrastructure/cloud/cloudfree>.
- [4] Open Access: Opportunities and Challenges [online] http://ec.europa.eu/research/science-society/document_library/pdf_06/open-access-handbook_en.pdf
- [5] OpenAIR – Open Access Infrastructure for Research in Europe.[Online] <http://www.openaire.eu/>
- [6] www.sas.com/en_us/in_sights/big-data/what-is-bigdata.htm
- [7] CloudCom2012poster
- [8] www.wikipedia.com