

*Real Time Data Integration Architecture for Hydrocarbon Exploration
and Production Data Consolidation*

By

KINGSHUK SRIVASTAVA

**COLLEGE OF ENGINEERING STUDIES-CENTER FOR INFORMATION
TECHNOLOGY**

Under the Guidance of

Dr. ASHISH BHARDWAJ, CHIEF INFORMATION OFFICER, UPES

Submitted



**IN PARTIAL FULFILLMENT OF THE REQUIREMENT OF THE DEGREE OF DOCTOR OF
PHILOSOPHY**

TO

UNIVERSITY OF PETROLEUM AND ENERGY STUDIES

DEHRADUN

January, 2013

©University of Petroleum and Energy Studies, Dehradun, 2012
All Rights Reserved.

DEDICATED
TO
My PARENTS

ACKNOWLEDGEMENTS

I wish to express my earnest gratitude to **Dr. Ashish Bharadwaj**, Chief Information Officer, University of Petroleum and Energy Studies (UPES) Dehradun, who has supervised my research work with a keen interest. His ever helping attitude, excellent leadership and dynamic personality have been a constant source of encouragement for me. And also, it is a great appreciation to him that he has provided a lot of facilities for all the experiments I am able to conduct in the different laboratories as well as undergo training in “Informatica & Cognos”.

I wish to express my sincere indebtedness to **Dr. Ashutosh Pasricha**, Manager, Oilfield Services and Real Time Solutions, Schlumberger (India), who has mentored me in my research work with zealously. His perpetually support all-embracing awareness in the field of hydrocarbon exploration and production and eternal espousal have been a prominent part of my research work. I would like to extend my earnest indebtedness for the training he enabled me to receive in “*Ocean Platform- Petrel*” and the actual understanding of research work.

I wish to express my gratitude to **Dr. B. P. Pandey**, Dean Emeritus, College Of Engineering, UPES, who facilitated me in finding the core problem area for my research work. He has always been a father figure and guided me through the ups and downs of my research days.

I would like to extend my honest gratefulness to **Dr. S. J. Chopra**, Chancellor, UPES, whose ever willing support to my research have been instrumental in completion of my work.

I wish to convey my sincere gratitude to **Dr. Parag Dewan**, Vice Chancellor, UPES, who had confidence in my abilities and gave me the chance to do research in UPES.

I would like to convey my earnest gratitude towards **Dr. Manish Prateek**, Professor & Head Dept. of CIT. UPES, for his constant understanding and support to my work.

My special thanks are to my Parents for their blessings and encouragement during the period of the study.

I am very thankful to my wife Monica for her excellent cooperation and support during the entire period of this research.

I wish to express my thanks to my sister Papia and my niece Titly for their blessings and encouragement throughout the period of my study.

My special thanks are also due to all the faculty members of Center for Information Technology for their help and support during the course of my study.

I wish to extend my thanks to my dear friend and colleague Dr. P. S. V. S Sridhar for his perpetual support and encouragement during the entire period of my work.

I want to extend my sincere thanks to Dr. Mousumi Dasgupta, National Training Manager at Indian School of Petroleum & Energy, New Delhi Area, India, for her belief in my abilities and her unwavering support for me.

Dehradun

January, 2013

(KINGSHUK SRIVASTAVA)

Declaration

"I hereby declare that this submission is my own work and that, to the best of my knowledge and belief, it contains no material previously published or written by another person nor material which has been accepted for the award of any other degree or diploma of the university or other institute of higher learning, except where due acknowledgment has been made in the text.

KINGSHUK SRIVASTAVA

Date:-

Place:- Dehradun

THESIS COMPLETION CERTIFICATE

This is to certify that the thesis on “*Real Time Data Integration Architecture for Hydrocarbon Exploration and Production Data Consolidation*” by **KINGSHUK SRIVASTAVA** in Partial completion of the requirements for the award of the Degree of Doctor of Philosophy (Engineering) is an original work carried out by him under our joint supervision and guidance.

It is certified that the work has not been submitted anywhere else for the award of any other diploma or degree of this or any other University.

External Guide

Dr. Ashutosh Pasricha,

Manager, Oilfield Services and Real Time Solutions, Schlumberger (India)

Internal Guide

Dr. Ashish Bharadwaj,

Chief Information Officer, University of Petroleum and Energy Studies (UPES) Dehradun

TABLE OF CONTENTS

	Page No.
ACKNOWLEDGEMENTS	i
DECLARATION	iv
THESIS COMPLETION CERTIFICATE	v
TABLE OF CONTENTS	vi
EXECUTIVE SUMMARY	xi
LIST OF ABBREVIATIONS	xvi
LIST OF FIGURES	xix
LIST OF TABLES	xxii
CHAPTER 1 INTRODUCTION	1
1.1. INTRODUCTION	1
1.2. MOTIVATION OF WORK.	2
1.2.1. REVIEW OF LITERATURE.	3
1.3.OBJECTIVE	5
1.4.METHODOLOGY	6

1.5. STRUCTURE OF THE THESIS	7
CHAPTER 2. REVIEW OF LITERATURE	8
2.1 ETL (Extract, Transform, Load)	8
2.1.1. INTRODUCTION	8
2.1.2. HISTORY OF ETL	9
2.1.3. ETL FRAMEWORK	13
2.1.3.1. HAND-CODED ETL PROCESS	14
2.1.3.2. TOOL-BASED ETL PROCESS	16
2.1.3.3. FIRST GENERATION – CODE GENERATORS	17
2.1.3.4. SECOND GENERATION – ETL ENGINES	19
2.1.3.5.ETL TOOLS TODAY	23
2.1.4. FUTURE TRENDS & REQUIREMENTS IN ETL TOOLS	25
2.2. SERVICE INTEGRATION APPROACHES	26
2.2.1. DATA AS A SERVICE.	26
2.2.2. INTEGRATION SYSTEMS	27

2.3. CLOUD COMPUTING	31
2.3.1. INTRODUCTION	31
2.3.2.LEXICON	34
2.3.3. ENVIRONMENTS OF THE CLOUD.	37
2.3.4. DEPLOYMENT TYPES	42
2.3.5 . CLOUD ENVIRONMENT ROLES	44
2.3.6. SPECIFIC CHARACTERISTICS /	47
CAPABILITIES OF CLOUDS	
2.4. REVIEW ON DATA INTEGRATION	52
2.4.1.INTRODUCTION	52
2.4.2.CHALLENGES OF DATA INTEGRATION:	55
2.4.3. APPROACHES TO INTEGRATION	59
2.5. INTRODUCTION TO EXPLORATION &	67
PRODUCTION DATA ORGANIZATION	
2.5.1. BASIC DATA INTEGRATION	67
CHALLENGES IN E&P SECTOR.	
2.5.2. TYPES OF DATA INVOLVED IN E&P	76
SECTOR.	

2.6. REAL-TIME DATA INTEGRATION.	86
2.6.1. INTRODUCTION	86
2.6.2. CURRENT WORKFLOWS AND PROCESSES INTEGRATION	88
CHAPTER 3. REAL TIME DATA INTEGRATION ARCHITECTURE (RTDIA)	89
3.1. REAL TIME DATA INTEGRATION ARCHITECTURE.	8.9
3.2. THE LOWER LEVEL ARCHITECTURE (EXPANDED).	95
3.2.1. DATA INTEGRATION LAYER	97
3.2.2.. HANA DATABASE	98
3.3. IDENTIFICATION & INTEGRATION ALGORITHM	101
3.3.1. INTRODUCTION	101
3.3.2. DATA MODEL	101
3.3.3. STORAGE SYSTEM	102
3.4. DATA IDENTIFICATION AND ASSOCIATION ALGORITHM.	103

3.5. WORKING OF THE ALGORITHM.	109
3.6. IMPLEMENTATION OF THE PROPOSED ARCHITECTURE.	109
3.6.1. THE DATA IDENTIFICATION SCENARIO.	110
3.6.2. SOFTWARE WORKING DESCRIPTION.	112
3.7. Testing.	112
3.7.1. TESTING ARCHITECTURE:	112
3.7.2. THE SETUP.	115
3.7.3. TABLE STRUCTURE:	117
3.7.4. THE FRONT END	124
3.7.4.1. ASSOCIATION RULES.	126
3.7.4.2. QUERIES FOR ASSOCIATION	127
3.7.5. WORKING OF THE SOFTWARE.	128
3.7.6. THE BENEFITS OF THE SYSTEM.	128
3.8. RESULTS	129
CHAPTER 4. CONCLUSION.	130
SCOPE OF FUTURE WORK	133

BIBLIOGRAPHY	135
BRIEF BIO-DATA OF THE AUTHOR	160

EXECUTIVE SUMMARY

The exploration and production sector of hydrocarbon is not only highly competitive but also requires an enormous amount of investment. The world market today is facing shortage of oil and gas supplies leading to constant increase of prices, thus the requirement for further new discoveries as well as prompt decision making has become more important than ever. Most of the old oil fields are depleting and the new reserves are in difficult and sometimes almost in inaccessible location for easy exploration and production. Improper or untimely decision can lead to huge financial losses and may cause irreparable damage to the formation. For maintaining a competitive edge in this ruthless market, accurate and efficient information and knowledge is required in real time or in near real-time environment. For extraction from existing reservoirs high speed and intelligent simulation solution becomes a critical factor. Under this highly complicated scenario Exploration and Production (E&P) companies requires to utilize and implement information technologies to streamline their processes as well as for efficient, accurate timely decision making.

Separate independent studies by International oil companies (IOC) corroborate the overall finding that “engineers spend half their time unproductively chasing data.” One IOC has identified standardization of data and information management practices as a foundation for enterprise wide work-process integration and automation. The company is developing production surveillance and optimization (PS&O) solution, which was justified on its

potential to minimize the amount of time required by engineers and technicians to access data and prepare it for analysis. Based on an internal survey of current PS&O processes across all production units around the world, the IOC identified that engineers and technicians spend an average of 44% of their time accessing information and preparing it for analysis.

The survey further revealed that currently a meager 9% of the operator respondents are able to get data automatically from their real-time systems into their engineering or geo-science analysis routines. The other 91% of operator respondents had to consume more than 50% of their time and human resources to identify, format, and prepare the data for the analysis tools. A staggering 55% of the respondents had less than 25% of their professional time available for analysis, decision, and action. These statistics show that there is a significant potential workforce that is sequestered in the data commute, unavailable for value-adding activities that utilize their engineering and geo-science education and experience.

To be able to support speedy, intelligent analysis for optimum utilization of upstream assets, high performance computers as well as large data storage infrastructure is required. These are the companies which were the first to deploy clusters and grids. The concept of data warehousing was a comparatively recent technology which came into the market with the advent of more powerful tools like ETL, and high level programming languages.

But in E&P sector the amount of data to be handled is in the range of 100s of TB, which creates a challenge for the IT technologies being implemented. The

technology has to be lot more sophisticated and the approach quite different from the conventional technology available in the market today.

There are a lot of opportunity and space for development in this area. With the evolution of Information Technology (IT) there are a lot of technologies available today which could be incorporated for the streamlining of the information identification and integration.

The biggest problem of the E&P sector is not only the amount of data it has to handle but also the diverse computer system platforms which are used. The IT infrastructure of the E&P sector normally contains an assortment of operating systems ranging from windows, Linux as well as may contain Macintosh. The jumble does not end here; diversity is present in databases maintained by an E&P company. Data could reside in a proprietary database like Oracle, MS-SQL server as well as open-source databases like MySQL. This heterogeneity in case of Data Base Management System (DBMS) as well as operating system is a problem when communication is required between them for any kind of data transfer. There are solutions present in the market like Open Database Connectivity (ODBC), or VM Ware (Virtual Machine) to handle the above stated problem of communication but each one of it has its own pros and cons. An exhaustive study on organization of E&P data revealed further snags in its total data organization architecture leading to major predicament in data identification and process of integration. The whole architecture of IT infrastructure is in a chaotic state at a very best of time.

If the architecture of IT infrastructure could be streamlined appropriately most of the complications could be resolved. To solve the problem of

communication between all the assortments of databases is to put them onto a common platform like a cloud structure for convenient communication between them. The proposed cloud platform described in this thesis Eucalyptus Enterprise Edition (EE) 2.0 is an open-source, Linux-based software architecture which has the ability to realize a scalable, proficiency-augmenting private as well as hybrid clouds, utilizing the available IT infrastructure. It being an Infrastructure as a Service (IaaS) provider the company can utilize and implement its own assortment of resources (hardware, storage and network) all over its distributed assets. The cloud could be deployed on the On-Premise Data Center and can be accessed over the companies' VPN (Virtual Private Network). As a consequence of which it enables a higher security to the confidential data of an E&P company.

The next step towards homogenization of the architecture is paramount for proper linking and removal of disorganized components in the IT Infrastructure. The propose architecture provides a perfect organized structure in the current scenario of an E&P company. The architecture streamlines all the different components of databases for the proposed real-time data integration process. Once the reorganization is done it becomes simple for an efficient algorithm for the extraction and integration of data from the source data marts to the target data warehouse.

The algorithm developed and described in this thesis looks into the different databases at located in different asset locations in a round robin basis within a scheduled interval of time. To be able to extract and integrate data at high speed for enabling real time data integration the algorithm needs to be simple,

light as well as agile. The algorithm designed here is very simple in its approach thus it is light. It generates a search command at a regular interval of time and follows all the enlisted distributed databases each at a time, whenever it encounters a new entry in it, it copies that entry into the central data warehouse with a time and date stamp. The next time the query visit the same database it looks for the last update date and time stamp and any entry later than that is again updated.

This enables a continuous real time data extraction and integration of E&P data into the central data warehouse enabling a faster and efficient process of Business analysis.

Through this process the basic problem identified in the beginning of the research of manually locating and extracting data for different analysis would become much more convenient and highly accurate.

List of Abbreviations

Abbreviation	Meaning
APIs	Application Programming Interface
BCP	Bulk Copy Program
BI	Business Intelligence
CAPEX	Capital Expenditure
COBOL	Common Business Oriented Language
CPU	Central Processing Unit
CRM	Customer Relationship Management
DBMS	Database Management System
DCS	Distributed Control System
DTD	Document Type Definition
DWH	Data Warehouse
E&P	Exploration and Production
EC2	Elastic Compute Cloud
EIP	Enterprise Integration Pattern
EME	Enterprise Meta Environment
ETL	Extract Transform Load
FDBMS	Federated Database Management System
G&G	Geological And Geophysical
GAV	Global As View
GIS	Geographic Information System

HPC	High Performance Computing
IaaS	Information As A Service
ILM	Information Lifecycle Management
IOC	International Oil Companies
IT	Information Technology
LAV	Local As View
MDM	Master Data Management
MSL	Mediator Specification Language
O/S	Operating System
ODBC	Open Database Connectivity
OEM	Object Exchange Model
OLTP	Online Transaction Processing
OPEX	Operating Expense
P2P	Peer to Peer
PaaS	Platform as a Service
PL/SQL	Procedural Language/Structured Query Language
PS&O	Production Surveillance And Optimization
QoS	Quality Of Service
RAID	Redundant Array Of Independent/Inexpensive Disk
ROI	Return On Investment
RTDIA	Real Time Data Integration Architecture
S3	Simple Storage Service
SaaS	Software as a Service
SCADA	Supervisory Control And Data Acquisition

SOA	Service Oriented Architecture
SQL	Structured Query Language
SQS	Simple Queue Service
TB	Terabyte
VM	Virtual Machines
VPN	Virtual Private Network
XML	Extensible Markup Language

LIST OF FIGURES

Number	Title	Page No.
2.1	ETL Processing Framework	12
2.2	Code Generators	18
2.3	ETL Engine	20
2.4	Conceptual View Of Cloud Computing	33
2.5	Non-Exhaustive View On The Main Aspects Forming A Cloud System	36
2.6	Modern Application Platform View	40
2.7	Data Integration Chart	54
2.8	Manual Data Integration Process	66
2.9	Upstream Domain Classification	69
2.10	Area Of Work	71
2.11	Rock Cycle	77
2.12	Seismic Survey Image	80
2.13	Well Logging Data	82

2.14	Reservoir Simulation Image	84
2.15	Current Organization Of Data	87
3.1	Architecture I	92
3.2	Expansion Of “Lower Layer” Architecture	94
3.3.	Basic Data Flow Diagram Of HANA	100
3.4	Component Diagram Of HANA Engine	100
3.5	Structure Of Data Model-1	104
3.6	Structure Of Data Model-2	104
3.7	Data Relationship	111
3.8	Architecture Of The Setup.	114
3.9	Primary Table Structure.	116
3.10	Data Cube Structure.	122
3.11	The Star Schema	122
3.12	User Interface -1	124
3.13	User Interface -1	125

List of Tables

Number	Title	Page No.
2.1	Various Generations Of ETL	10
2.2	List Of Programs For Development Of ETL Tool.	15
2.3	Advantages & Disadvantages Of Hand Coded ETL	16
2.4	Advantages & Disadvantages Of Code Generator ETL	19
2.5.	Advantages & Disadvantages Second Generation ETL	22
2.6	Advantages & Disadvantages Present ETL Tools.	24
3.1	Operational Data Types	72
3.2	Operations For Exploration & Production	73
3.3	Asset	117
3.4	Well	118

3.5	Assay	119
3.6	Structure	120
3.7	Steam Injection	120
3.8	Production	121
3.9	Temperature & Pressure	121

CHAPTER 1

1.1 INTRODUCTION

The exploration and production sector of hydrocarbon is not only highly competitive but also requires an enormous amount of investment (according to reports by HIS-[Information Handling Services] the capital expenditure was \$641 billion in 2012) [154]. The world market today is facing shortage of oil and gas supplies leading to constant increase of prices, thus the requirement for further new discoveries as well as prompt decision making has become more important than ever. Most of the old oil fields are depleting and the new reserves are in difficult and sometimes almost in inaccessible location for easy exploration and production. Improper or untimely decision can lead to huge financial losses and may cause irreparable damage to the formation. For maintaining a competitive edge in this cut throat market accurate and efficient information and knowledge is required in real time or near real-time environment. For extraction from existing reservoirs high speed and intelligent simulation solution becomes a critical factor. Under this highly complicated scenario Exploration and Production (E&P) companies requires to utilize and implement information technologies to streamline their processes as well as for efficient, accurate timely decision making.

To be able to support speedy, intelligent analysis for optimum utilization of upstream assets, high performance computers as well as large data storage infrastructure is required. Oil companies are the companies which were the first to deploy clusters and grids. The concept of data warehousing was a comparatively recent technology which came into the market with the advent of more powerful tools like “ETL” (Extract. Transform. Load), and high level programming languages makes it possible today.

But in E&P sector the amount of data to be handled is in the range of 100s of “TB” (Terabytes) [153], which creates a challenge for the IT technologies being implemented. The technology has to be lot more sophisticated and the approach quite different from the conventional technology available in the market today.

1.2. Motivation of Work.

Separate independent studies [36] by “International oil companies” (IOC) corroborate the overall finding that “engineers spend half their time unproductively chasing data.” Consistency of data and information management customs has been acknowledged by one IOC as the basis of integration and automation of the whole enterprise [36]. It has ventured into building a software package called “production surveillance and optimization” (PS&O) solution, and under the circumstances it has a high rationality to go for it to reduce the unproductive time spent by engineers just preparing data for analysis. According to an internal survey conducted the company a staggering 44% of time is spent by

a technician or an engineer just for accessing and processing data for analysis purpose [36].

Further assessment of the same survey discovered that paltry 9% [36] of the operators were capable of retrieving data from the real time systems automatically, to feed into their geo-science or engineering routines. Whereas, 91% [36] of the operators conveyed that they on an average spend 50% of their time in identification and preprocessing of data, to enable it for different analysis tools. The results further showed that the time available for analysis and decision making was less than 25% for 55% [36] of the respondents. The above mentioned survey pointed out that a substantial amount of workforce is caught up in unproductive activities rather than utilizing their highly skilled potentials in productive and ground breaking work [36].

1.2.1.Review of Literature.

The primary literature which states the above given problem was published in SPE journal September 2009[36]. Which gave a clear identification of the problem in E&P data integration, identification and analysis? This states that almost 44% of productive working time is lost in unproductive and repetitive work. If this problem could be circumvented with the present technologies a huge amount of revenue (tangible & intangible) could be saved.[36]

The primary application of ETL techniques for the purpose of solving the problem of identification and cleaning of data could be easily used. Literature contains

many ideas put forward by many researchers to this end. “*A conceptual modeling of ETL processes*[2]” was published by Alkis Simitsis[2] and his partner. Another paper *Spatial ETL Tools: The Bottom Line of an Enterprise GIS*[155] was discussed by S.Raghavendran. It discussed the ETL tools which are being used for the purpose of data integration in GIS.

The ETL taxonomy was discussed in a paper written by Panos Vassiliadis et.al.[101]. In this paper the writer and his associates presented the method for identification of generic properties that characterize ETL activities. The concept of this paper is identification of generic properties that characterize ETL activities. For it they follow a black-box approach and provide a taxonomy that characterizes ETL activities in terms of the relationship of their input to their output and provide a normal form that is based on interpreted semantics for the black box activities. The proposed taxonomy can be used in the construction of larger modules, i.e., ETL archetype patterns, which can be used for the composition and optimization of ETL workflows.

The next important way of handling the problem is the use of Information Lifecycle Management which was nicely put forward by Mr.Paresh Chatterjee [29] in PCQUEST February. In this paper the “Storage Tiering” is discussed. Here, the capacity to be provisioned is divided into separate pools of storage space with various cost/ performance attributes. At the top resides the Tier 1 pool, which is the most expensive but high performing nonetheless. The bottom tier is occupied by more cost-effective storage arrays. The concept to devise a

sophisticated software layer that intelligently places data into the different tiers according to their value is also discussed. This concept is variously known as data classification or “Information Lifecycle Management” (ILM) which is highly required in my research work.

Another enlightening article was published in PCQUEST February edition 2010 by Rahul Sah, in which the author discussed the “advanced analytics”.[103]. Advanced analytics incorporates methodologies for answering future-oriented, proactive and predictive questions, as well as streaming data; so that consequent business decision is based on real-time questions about what's going on now. Advanced analytics leverages the same core features of typical analytics solution, i.e. reporting, dashboards, visualizations etc. but takes the analytics power to several steps further. Such powerful analytics capabilities would be required by enterprises where they generate bulk of data and there is a need to analyze that in real-time.

The above stated preliminary review of literature hints at the different technologies which could prove to be the avenue ahead to solve the problem at hand. The ETL technology could be the key to the solution of the problem if it could be modified and manipulated according to the specific requirement. This basically leads to the identification of methodology of the research work. Further in the second chapter a thorough review work of few related technologies is given in detail.

1.3. Objective.

The objective of this work would be as follows:

- To define and design an effective and efficient architecture.
- To develop an intelligent algorithm.

To facilitate a seamless real time data integration environment with ultimate aim of reducing the amount of time lost in information handling by E&P workforce in searching for data, integrating it from multiple sources, and preparing it for analysis in applications.

1.4. Methodology

- 1) Identify the different types of data generated in the following operations.(see chapter 2.6.2 for details)
 - i. Geological Data
 - ii. Seismic Data/geophysical data.
 - iii. Drilling Data
 - iv. Well Log Data
 - v. Production Data.
- 2) Specify the components for the required architecture.
- 3) Design efficient data integration architecture.
- 4) Develop the required data integration algorithm.
- 5) Identify the areas of further improvement.
- 6) Give suggestions for future improvement.

1.5. Structure of Thesis

The thesis is divided into four chapters. The first chapter discusses the problem and some review literature which authenticates the problem existing in the market. It also looks into the present technologies and its pros and cons for the specific market of Exploration and Production sector.

The second chapter reviews the different technological researches and development which has been done till date concerned towards the problem at hand. It also looks into the different aspects of Exploitation and Production Data Organization in attendance.

The third chapter contains the design and the components of the required architecture for the real time data integration in the E&P sector with the concerned discussions and arguments. The later part of the chapter states and explains the algorithm and the software developed for the testing of the algorithm.

The fourth chapter is the conclusion of the thesis, which looks into the problem; the solution formulated its limitations and scope of future work.

CHAPTER 2

REVIEW OF LITERATURE

This chapter entails basically three technologies namely ETL (Extract, Transform, load), Cloud Computing, data & service integration techniques and also the current scenario of the Exploration and Production Information Technology architecture. Here an attempt has been made to comprehensively study the different technologies available today in the market which could be utilized and implemented for the creation of an efficient architecture.

2.1 ETL (Extract, Transform, Load)

2.1.1. INTRODUCTION.

The abbreviation ETL is expanded as Extract, Transform, and Load, essentially it explains a process where data is episodically extracted out of the source systems, transformation of data into a coherent format (Data Cleaning), and loading of data onto the target database normally referred to as Data Warehouse. In the course of ETL process, extraction of data occurs from an On Line Transaction Processing (OLTP) database / non-OLTP database system, transformation of the same extracted data essentially meaning conversion of heterogeneous data types into a consistent format and equate with the data warehouse schema. Ultimately the data is uploaded and saved in the target data warehouse database known as loading.

Unpretentiously it is replication of data into different databases with various alteration processes in between. Basically ETL for data warehousing should be regarded as a process rather than a physical implementation.

2.1.2. HISTORY OF ETL.

With the gradual evolution of Data Warehousing, organizations required a process to load and maintain data in a Warehouse. ETL process evolved and gradually took control over the Data Warehousing market to fulfill this requirement. Initially, organizations developed their own custom codes to perform the ETL activity which was referred as Hand-coded ETL process[107]. Since the process was lengthy and quite difficult to maintain, vendors started developing of the shelf tools which could perform the same task as that of Hand-coded ETL but in efficient manner. In this era of ETL tools, the market saw different generations, different types and tools with their own pros and cons.

Table 1 gives the brief overview of ETL history starting from hand-coded ETL tools to ETL tools available today in the market.

ERA	TITLE	SIGNIFICANCE
Early 1990	Hand Coded	Custom Codes (Hand written)
1993-1997	1 st Generation Tools	Code Based Tools
1999-2001	2 nd Generation Tools	Engine Based Tools
2003- 2006	3 rd Generation Tools	Efficient Tools
2007- 2011	Parallel ETL Processing Tools	Intelligent Search & Optimize
2011- Till Date	In Memory Computing (HANA)	High speed processing, and handling of huge data sets.

Table 2.1: Various Generations of ETL[107]

The generations tabulated above are discussed below in details.

The business requirement should always be the principal concern for a developer before any kind of implementation of ETL tool. Not only the developer should concentrate on achieving the business requirement but also should be able to do it in an efficient way.

To illustrate the above situation, if a system requires loading 100 TB of data sets onto the target database per day, it should not only do it accurately but also be able to do it efficiently, like an exchange partition to increase the performance. The actual implementation process of ETL varies from data warehouse to data warehouse and even at departmental data marts within the same data warehouse.

Initially in the 1990's[88], all most every establishment developed and owned their personal custom made codes for pulling out data from operational system and extracting and transforming data from operational systems and insert it into data warehouses. In spite of the different way these systems are implemented the purpose of each of these are mutual. They essentially shift data from one database to another with some changes to the schema of source data. An ETL system consists of four distinct functional elements:

1. Extraction
2. Transformation
3. Loading
4. Meta data

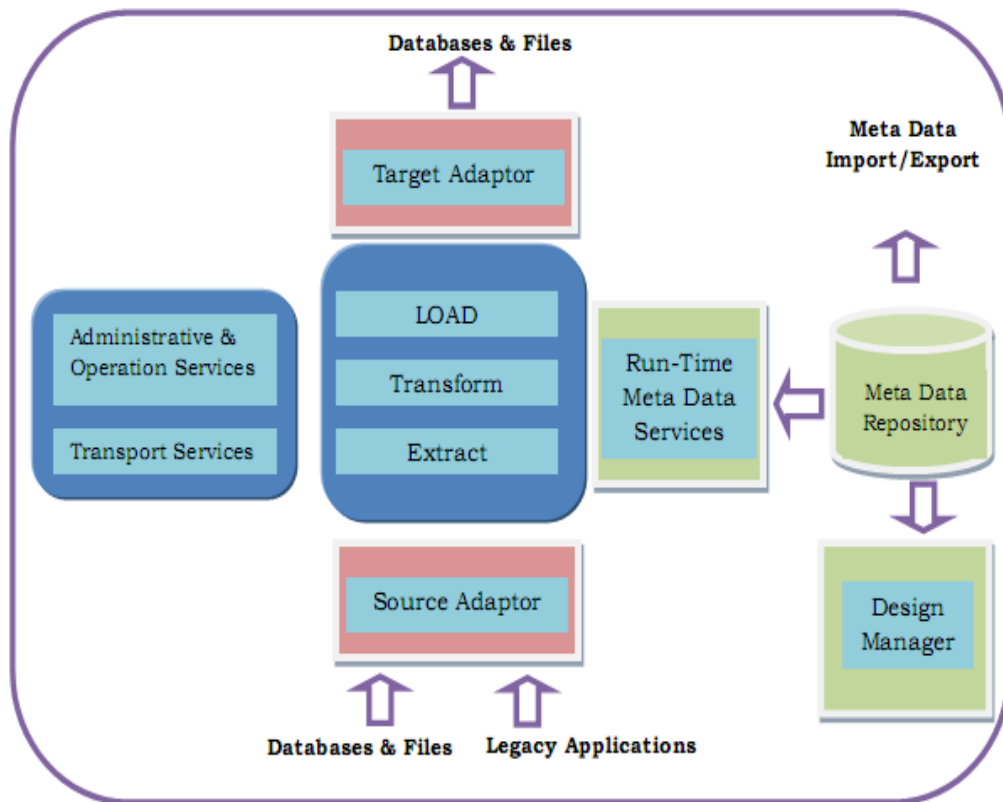


Figure 2.1: ETL Processing Framework [88]

2.1.3. ETL FRAMEWORK

The figure (Fig.2.1) depicts the basic components involved in ETL processing.

The bullets below describe each component:

- Extraction: This process of extracting or retrieving of data from source systems (which normally comprises of an OLTP, legacy, mainframe, flat files, etc or any combination of it) data using adapters, like native SQL formats, ODBC, or flat file extractors. These adapters consult metadata to determine which data is to be extracted and how [101].
- Transformation: This process transforms or converts the extracted data into a consistent data warehouse schema by applying the predefined rules onto it. This process is also responsible for validating, accuracy, type conversion of extracted data and business rules application. It is most complicated of the ETL elements [101].
- Load: The ETL Load element is responsible for loading the transformed data into data warehouse using target data adapters such as SQL loader, Bulk process, BCP, etc [101].
- Metadata: The metadata part of ETL system maintains all the required information about the data. The metadata repository makes metadata available to the ETL engine at run time. [101,107]
- Administration and Transport services: The ETL transport service consists of network and file transfer protocols to move data between source and target systems. The utilities enables administrators schedule, run, and

monitor ETL jobs as well as to log all events, manage errors, recover from failures, and reconcile outputs with source systems.

The above described components were used to be manually coded using native SQL codes, C, COBOL and other programming languages. Today, these components come with most vendor-supplied ETL tools in which all of these components and functions are combined together to create a single, integrated package.

The generations tabulated above are discussed below in details.

2.1.3.1. HAND-CODED ETL PROCESS

At the beginning developers used custom codes for the performance ETL operations. The programs which were written in this method were not only lengthy but also very difficult to document. The developer normally used different programming languages to perform the task. Normally it used to be a combination of technologies like SAS, Database, Perl, Shell, etc. [107]

PROGRAMS	APPLICATIONS
SHELL	Wrapper Scripts
PERL SCRIPTS	Data cleansing, pattern matching, auditing.
SAS	Reading source files and applying transformations.
ORACLE	Bulk loading & Direct loading.
AUTOSYS	Automation of the process.

Table 2.2: List of programs for development of ETL tool.[107,159]

These custom programs were not much of a viable option as there were a lot of issues and problems with this method of ETL process.

Advantages	Disadvantages
Manual creation of Meta Data gave direct control over the organization and running.	Continuous modification and rewriting of codes were required increasing the overall project cost.
Easy code testing for the availability of automated unit testing tools	Maintaining of separate metadata table was required and any changes required

Most Flexible & Customizable	<p>manually changing of the entire table.</p> <p>Single threaded and slow speed of execution.</p> <p>High development effort and difficult testing.</p>
------------------------------	---

Table 2.3. Advantages & Disadvantages of Hand Coded ETL

2.1.3.2. TOOL-BASED ETL PROCESS

To avoid the above essayed overheads caused by hand-coded ETL process, vendors started developing ETL tools to perform extraction, transformation and loading process. The most important aspect of these tools is it generates and maintains centralized metadata repository. With the development in computing capacity and distributed computing systems, and as business intelligence made its debut, the first ETL solutions were introduced. In the beginning ETL provided the ability to extract the data from mainframes and load into target database. Today the ETL tools have matured to provide user-friendly GUI's, performance benefits and additional functionalities.[115]

There is still a debate about whether ETL engines or code generators offers the best functionality and performance. Since that time, several generations of ETL have been produced.

2.1.3.3.. First Generation – Code Generators

To get rid of writing the complex hand-written codes, vendors started developing the ETL tools in mid-1990s and they started producing the legacy code generators.[136]

The code generating tools at that time were mostly based on COBOL as data was basically stored on mainframes. Programs for extraction was mostly written in batch mode which automatically generated source codes for compilation, scheduling and running. The data extraction from source files, transformation and loading of the data in database process used to run on server. These were single threaded programs which did not support parallel processing.

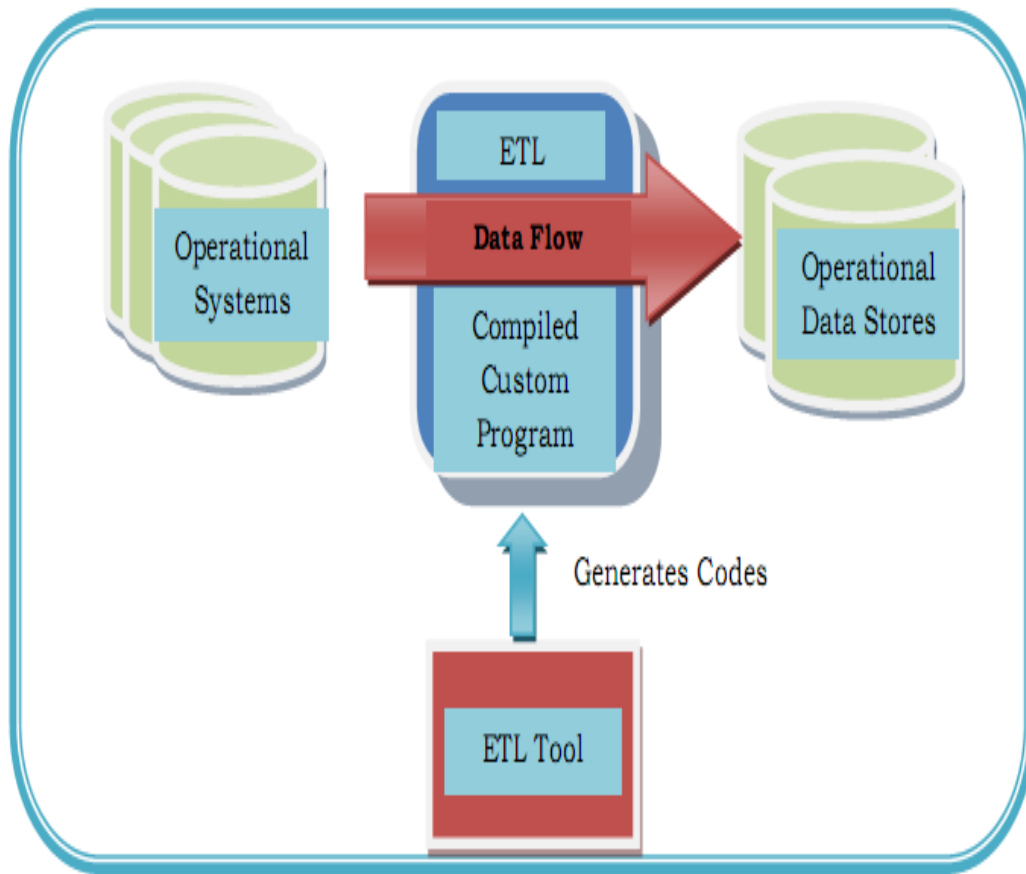


Figure 2.2: Code Generators.[136]

Advantages	Disadvantages
Very powerful in extraction of data from legacy system	Requirement of profound knowledge in COBOL or C.
Native compile code enhanced the performance.	Very low success rate in handling of large relational databases. Parallel processing not supported Manual coding required in many type of transformation.

Table 2.4. Advantages & Disadvantages of Code Generator ETL

SAS Institute Inc started releasing its SAS Version 6 series in 19156 which supported base SAS along with MACRO facility, ODBC support, data step debugger, etc and was running on Unix and Windows. In 1997, SAS warehouse Administrator was released. [136]

2.1.3.4. SECOND GENERATION – ETL ENGINES

In the late 1990’s vendors started delivering the “engine based” ETL tools to automate more of the ETL process.

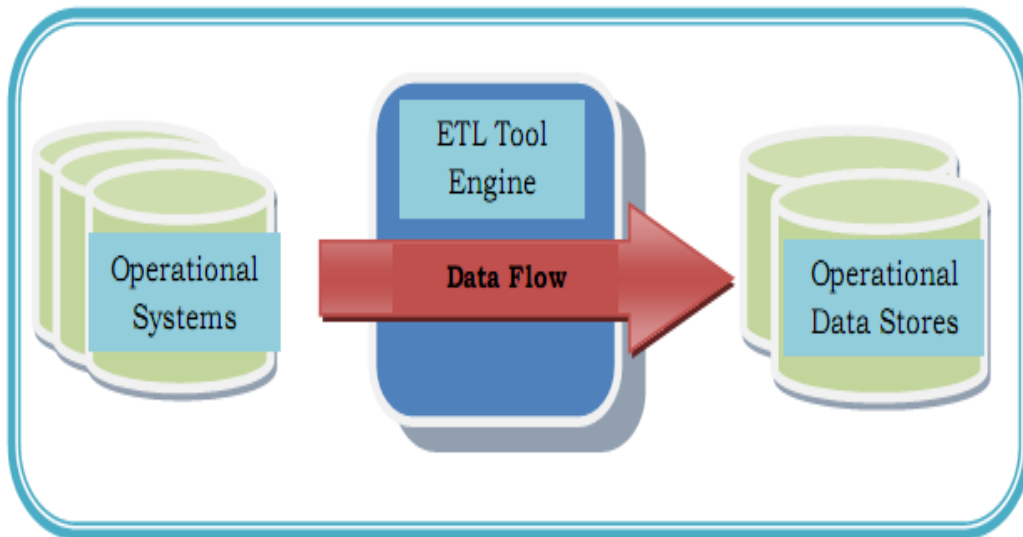


Figure 2.3: ETL Engine[136]

In the first generation tools codes were generated but in an engine based ETL manages the process through an internally present engine. The engine applies language interpreters for processing of ETL workflows at runtime. The predefined ETL workflows are stored in a Metadata repository enabling the engine to determine how to process incoming data at runtime. This enabled the user to have knowledge of just a single programming language which is the language of the ETL tool to solve all the problems. Another significant characteristic of an engine-based approach is that all processing takes place in the engine, not on source systems. Some of the engines support parallel processing, to achieve it partitioning needs to be done at the server manager. Also in engine-based tools metadata repository is maintained at different server which needs to be connected from client application.

Advantages	Disadvantages
Features like Graphical User Interface (GUI) and data transformation present	Extremely slow process of transformation when the data volume is high.
Increased efficiency as it is Multithreaded	Transformation process dogged by bottleneck problem
Substantially integrated and automated functions.	For better performance dedicated administrator for engine based tools.
Functions like aggregation, scheduling, transformation and monitoring are incorporated	

Table 2.5. Advantages & Disadvantages Second Generation ETL

❖ **Examples of Second Generation ETL Tools**

1. Power mart 4.5 developed by Informatica Corporation developed in 1999.
2. “Datastage” 4.0 developed by Informix corp. in 2000 which is now owned by IBM corp.

2.1.3.5. ETL TOOLS TODAY

Today's ETL tools are highly embedded with transformation features can support surrogate key generation, multi-dimensional designs, various transformation functions multiple input or output database or flat files, and native database or O/S utility. They have their own independent and internal metadata repositories which could be different from the data warehouse metadata repository. These are able to eliminate developing and maintaining of complex routines and transformations in ETL workflows. They also provide user friendly GUI features which enables the user not to have a very deep knowledge or training in the programming languages. They have features like bulk loading, monitoring, incremental aggregation, scheduling, etc. [115]

Previously only COBOL was used for the generation of code but today it supports many other platforms like SQL. Most of the code generating tools is using SQL to generate the code and also supporting all SQL features. The most important feature of code generator is its compile code which can run on any platform. The compiled code is much faster and distributes load across multiple platforms to increase the performance. The Engine based tools are also emerging very fast. They are coming with the ability to perform complex transformations at a faster speed of execution.

Advantages	Disadvantages
Rapid execution of Complex transformations handled effortlessly	Inability to support real time application
Supports heterogeneous platform for optimum performance	A number of ETL tools rely upon the type of source databases.
Sustains various category of Parallelism(Data, Component, Pipeline)	A majority of the contemporary tools don't support metadata level integration.
Ability to read XML (Extensible Markup language) data proficiently.	
Schedulers are incorporated to arrange workflows.	
Every Data Warehousing concepts holds good. E.g. slowly changing dimension, normalization, de-normalization, etc.	
Version controlling is present so that overwriting of the source code does not happen. E.g. Ab-initio's EME (Enterprise Meta Environment)	
Specific Debugger present.	

Table 2.6. Advantages & Disadvantages Present ETL Tools.

2.1.4. FUTURE TRENDS & REQUIREMENTS IN ETL TOOLS

Today we are looking into a market scenario where real time data is required for taking decisions, a decision which can make or break an organization, cost the company a huge loss or may cause fatal accident.

For this kind of high performance requirements, ETL tools needs to consists of parallel processing capabilities which can assembles and validate data against the existing and next coming data. The Data Warehouse in question needs to be highly active to serve up to the high requirement. [115]

If we are talking about the Exploration and Production (E&P) of Hydrocarbon sectors the major challenges in data warehousing is the enormous volume of data it requires to handle each day and is increasing in an exponential manner. Normally a data warehouse is loaded during off-peak hours and queries, reports and views are generated on normal business hours for decision making. In this concept whatever may be the volume of data we are looking into it has to be done in off-peak hours and cannot be extended.

By keeping the increasing data volume in mind and time constraint of loading hours, ETL has to be efficient enough to adjust with this day-by-day increasing data volume. For this to be possible today's ETL tools should be able to substantially shorten the delivery time and improve application maintenance. Also

the tools should possess intelligent applications which take into consideration the previous processes.

The developer sometimes find difficulties in identifying and removing bugs for the lack of proper error messages, so proper error messages should be incorporated into it for creating a tool which has a very low downtime.

One of the best ways to incorporate a high throughput is to have web based workflow solutions. Today the business world is looking into the cloud based architecture where XML data is the major share holder and for easier and swifter data integration organizations are looking forward to web based data sources.

With the above mentioned features the ETL tool would become an efficient and highly effective system which should provide for the future market needs.

2.2. SERVICE INTEGRATION APPROACHES

2.2.1. DATA AS A SERVICE.

Data as a service defines data and service together where data is considered as specific type of service. Truong and Dustdar, [132] was able to describe WSDL or REST full [105] interface. Under the “*Service Oriented Architectures (SOA)*” [60,145] data services are understood as value added services. Amazon Simple Storage Service (S3) [3] is a very good example of Data as a Service concept which provides unlimited storage through a simple service interface.

In this thesis Data as a Service approach is applied as data has to be prioritized according to the requirement algorithm which further adds value to it, conforming to the definition of Data as a Service.

2.2.2. INTEGRATION SYSTEMS

For over thirty years data integration has been an area of interest and research for researchers around the world specially people who are in the field of databases, semantic web, data warehousing, ETL, artificial intelligence etc. Different issues, limitations and solutions has been discussed, formulated and improvised in the field of data integration in these years. An enormous number of projects have been produced based on these conceptual approaches and models creating difficulties in classification of these works on the basis of comprehensive and shared criteria.

Among many the most discussed data integration architecture is “Mediator” [139]. It is basically built on the concept of creating a Mediated schema (or Global schema) for synthesizing the source structure which requires to be integrated. In this architecture the collected data managed in a common way gives the user a virtual global perception for his queries generated enabled by the median schema. By automatic “unfolding-rewriting operations” a single query generated by the user is converted into multiple sub queries and declared as a set with respect to the mediated and source schema. To reconstruct it back the output of these sub queries are processed by data reconciliation techniques. The

correspondence between the source and mediated schema data turns out to be the most important aspect of this architecture [1].

Global-as-view (GAV) and local-as-view (LAV) are the two basic concepts of data integration systems used for mapping purpose which are mentioned in the literature. In GAV the queries generated on the source is converted in relation to the contents of the mediated schema, whereas in case of LAV the data source is expressed as the view on the mediated schema.

The above mentioned concepts are amalgamated into a single model known as global-local-as-view (GLAV). In this model features and properties of both the above mentioned approaches are coupled together [50]. In GAV scenario the major issue occurs while updating the median schema, if the source changes it would have reverberating effect on the mediated schema and the supporting applications. In LAV for the creation of query manager a huge number of issues are present for query rewriting. Since then many systems has been designed and developed which implemented this architecture. (see [59] for a survey).

Among the initial data integration projects TSIMMIS [78] system was the first to follow a "structural" approach to create a GAV architecture for structure and semi-structured source data. In this project the concept of wrapper was implemented which translated the source data model into a self describing model OEM (Object Exchange Model), the MSL (Mediator Specification Language) which were predefined manually by the designer for source integration was used

to deploy these rules. Arbitrary or rather specifically recursive views are defined at the mediator layer with the help of MSL in TSIMMIS.

A source and query independent mediator was conceived in the *Information Manifold system* [50]. In this architecture the integrated schema as well as the description of the source was manually defined by the designer manually. An architecture based on the wrapper model was applied for describing local source data utilizing an object oriented language by Garlic [25]. In this concept data and transformation of schema were handled uniformly. A query processor is present which ultimately executes and optimizes different queries over heterogeneous data sources applied on an object extended SQL system was the main component of Garlic. The wrapper enables data transformations between the source and the middleware directly or by creating views and the system interacts directly with it. Garlic objects is an abstraction layer present in Garlic wrappers use an abstraction the layer which describes data from heterogeneous sources into an uniform format is known as Garlic Definition [25] Language. These are then utilized to define the global schema.

By encapsulating multiple data sources and utilizing these objects views are created. The problem of schematic heterogeneity is acknowledged by the author in this architecture but the author does not try to tackle it directly. To tackle this problem a semi-automatic tool was developed named “Cilo” supports the data transformation and integration processes.

This project “Cilo” [44] made a pioneering contribution in schema mapping, it conceived a tool which created mapping between source and target data schemas in a semi-automatic approach.

Another prevalent problem in data integration architecture is the size of the ontology which requires to be handled; there are limitations when large ontology is required to be handled.

Among the many proposed approaches most could be classified under two categories of approaches [100]: Several approaches have been proposed to extract modules from a given ontology, that can be mainly classified in two different classes of approaches [100,58,69,135] Description Logics semantics and [94,37 ,117] graph-based representations of ontology.

There are different software which are being offered by big IT majors for data integration like,

I. Oracle Data Integrator [136]

II. Microsoft SQL Server Integration Services[136]

These tools basically provide ETL capabilities for Data Ware housing. Among the open source community “Talend” [136] was one of the forerunners. Initiated by a French company, it is based on the “Eclipse”[136] platform for ETL tools. Among its main features it has over 400 connectors some of which are application specific, it also has data profiling and cleansing capabilities.

2.3. CLOUD COMPUTING

2.3.1. INTRODUCTION

Cloud computing is a computational model in which assorted, heterogeneous and mutually excluding systems amalgamated together over a coherent and consistent collection of different kinds of networks which could be infinitely scalable in form of infrastructure so that it can give unconstrained support to data computation, applications and storage. This technology has revolutionized the IT sector with unbelievable lowering of input cost by an organization for availing and utilizing application hosting, content storage and delivery of computational information.

Cloud computing has brought about a paradigm shift in the potential of the computational technology as a whole. Primarily the technology which used to be out of hand due to the cost factor has become totally reachable in many aspects. In theory the concept of “Reusability of IT capability” is made possible by this new concept. Till now the input investment for an infrastructure has been capitals intensive, which sometimes were beyond the resources of an organization. This technology enabled an environment where the concerned organization would have to pay as per their requirement.

This concept is not that new, it all started with the concept of virtualization in networking, storage or application. Cloud computing brings together the

traditional concepts of “Grid Computing” , ”Distributed Computing” , “Utility Computing” , “ Autonomic Computing”.

Forrester defines cloud computing as:

“A pool of abstracted, highly scalable, and managed compute infrastructure capable of hosting end-customer applications and billed by consumption.” [162]

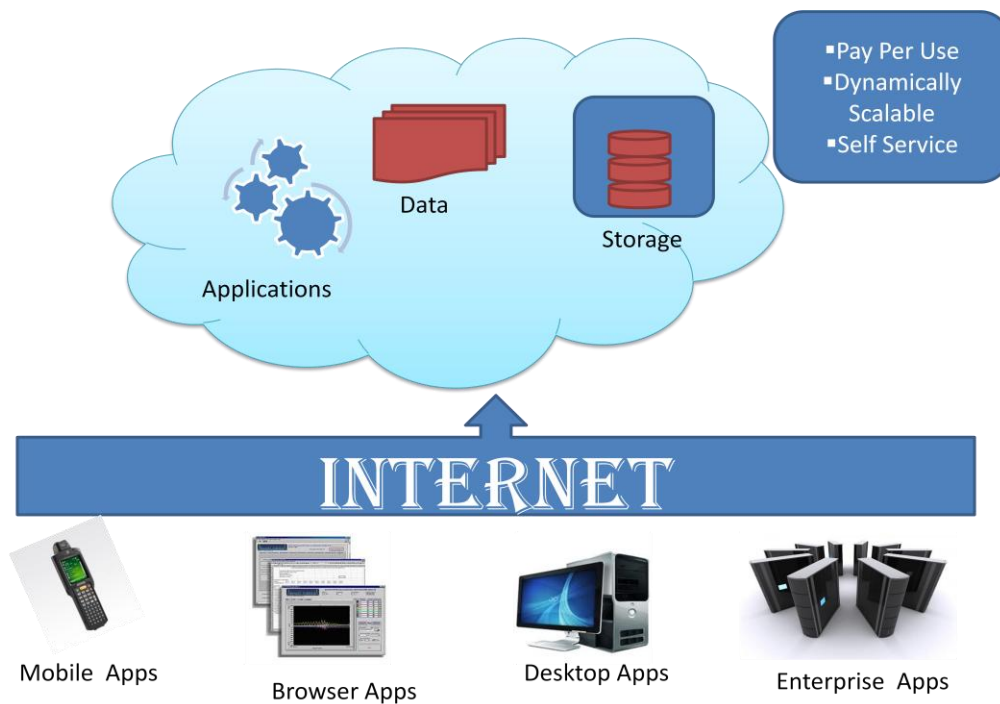


Figure2.4: Conceptual view of cloud computing Cloud Computing

Defining a cloud is as various as its applications, usage and the range of its implementation. Here no new definition or terminology is fashioned but an attempt has been made to identify and enumerate the different aspects and terminologies being used and implemented in the field of cloud computing.

2.3.2. LEXICON

In general a cloud could be described as a highly flexible implementation environment of diverse resources providing services on demand at a flexible rate according to the determined units of services and predefined quality. Basically it is a metered service provided to any customer globally at a very reasonable rate without any difficulty caused due to heterogeneity of operating system or any proprietary software applications.

Another way of describing cloud is that predominantly it is a platform which enables execution in diverse forms across multiple resources as well as in many cases across different enterprises. There are many distinctions between cloud types but the common criteria which binds them is that all of them are able to enhance resources and services either directly or indirectly. Another common feature among these clouds is that all of them augment manageability, platform independence and elasticity of a system.

It is assumed that some features like precise level of quality, performance criteria's as well as energy consumptions would be taken into consideration for

while allowing resource integration irrespective of any boundaries due to organization or stakeholders.

Manageable word signifies automatic insurance of predefined quality parameters, to be more precise, cloud can be explained as a compilation of diverse infrastructure/platform which facilitates implementation of services or applications in a more controllable and flexible mode. Controllable signifies consistency to the pre-identified and delineated quality measures inevitably, and flexible implies unrestricted scalability as per the user requirement of data and/or resources. The design of a cloud would actually vary significantly according to the performance and employment. The characteristics are defined by the capacity and facility of the cloud system according to the situation.

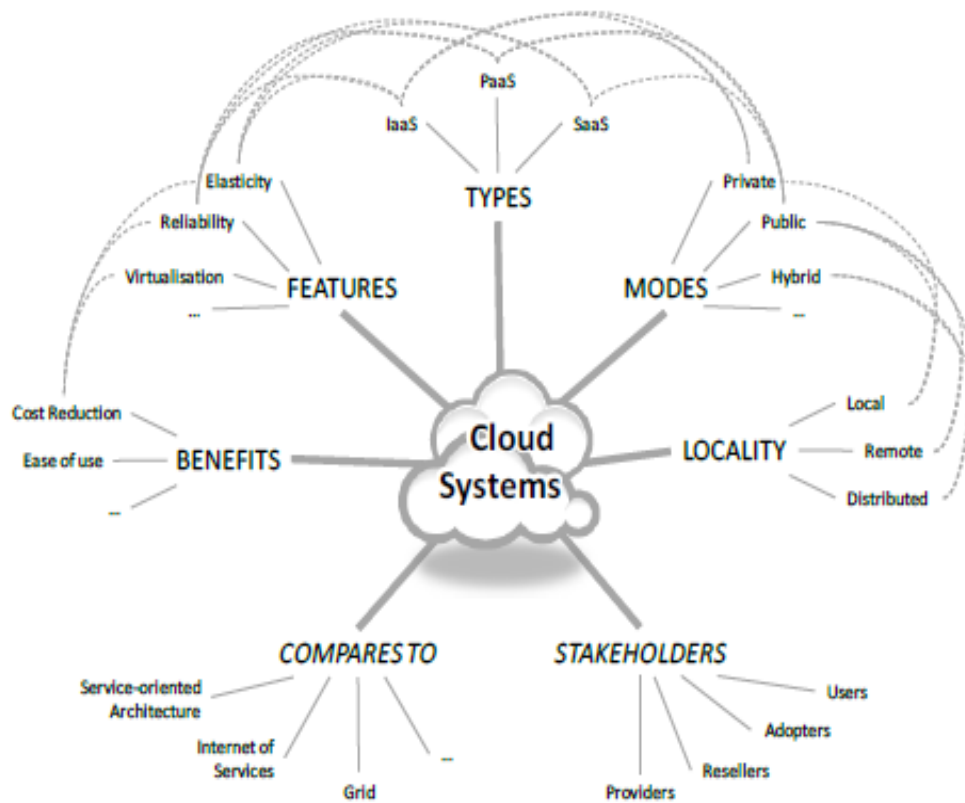


FIGURE 2.5: Non-Exhaustive View On The Main Aspects Forming A Cloud System[145]

2.3.3. ENVIRONMENTS OF THE CLOUD.

Typically Cloud exposes a particular functionality in accordance to the precise mode. Generally each cloud is relevant to a precise facility resembling Software Infrastructure, or Platform, but there are no limitations on the diverse type of facilities offered. Normally PaaS (Platform as a Service) providers tender specific applications as for example Google docs are included with Google App Engine. Therefore usually they are called “components” [53].

In terms of terminology the different literatures has slight variations among themselves due to the overlapping of various applications making it intricate for exact differentiation. In some literature a lot of popular terms are introduced which are far from technical terminology, making the terming of cloud more varied

The subsequent list defines the prominent types of clouds :

I. Infrastructure as a Service (IaaS) / Resource Clouds.

Through service interfaces it provides superior virtualization capabilities and renders highly managed and scalable resources as a service to a user. In other terms it makes available different services to a user imperative of the hardware or software constraints of the user. [53]

II. Compute Clouds:

In this model raw access to hardware computational power like CPUs (Central Processing Unit) is provided to a user, it thus tenders the facility of computing resources. Previously such an access to the computational hardware was not available by itself, but came as a part of “virtualized environment”. It could sometimes be confused with PaaS (Platform as a Service) like hypervisors which actually provides full software stacks for development and deployment of applications and are characteristically virtualized for execution of cloud services and applications, again in case of IaaS (Infrastructure as a Service) provides extra facilities in excess of just simple computation. Examples: Elastichosts, Amazon EC2, Zimory.

III. Data & Storage Clouds:

It is basically an online highly flexible comparatively cheap and one of the most secured way of storing data today. It bestows reliability and security to dynamic sizes of data and is charged accordingly. There are now more than one option available in the market today in this sector other than big players like Amazon S3, SQL Azure, Google Drive, there are companies like Dropbox, Spideroak etc. [53,4]

IV. Software as a Service (SaaS):

This type of cloud is also acknowledged as *Service or Application Clouds*.

Basically it pertains to discharge of precise business functions and processes which are made available through cloud potential. Through cloud infrastructure or platform applications or services are provided to the user in this type of cloud.

Normally it offers standard software applications or their functionality. [4]

Examples: Google Docs, Salesforce CRM, SAP Business by Design.

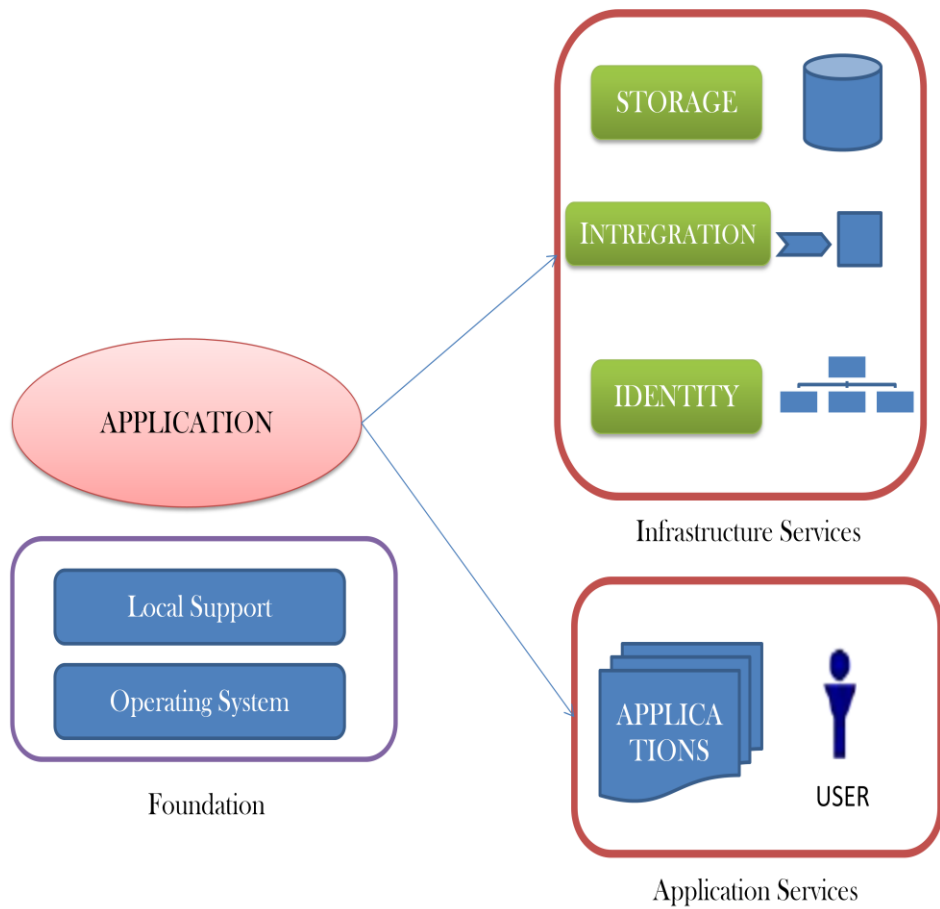


Figure 2.6: Modern application platform view

V. Platform as a Service (PaaS).

In this model the cloud creates a platform with whatever in terms of software and hardware resources is required to develop and host a certain application at a convenient pay per use basis to the user. In other words computational resources are made available to the user for their personal use. These systems contains dedicated APIs (Application Programming Interface) which enables and controls the dedicated server engine controlling and executing the user demands like access rate etc. On the downside of this model every service provider has their own design of APIs in accordance with their strategic capabilities, this causes portability issues if a user wants to migrate from one service provider to another. Efforts are being made to bridge this gap by introducing standard programming models in cloud, like MS Azure. [4]

Examples: Force.com, Google App Engine, Windows Azure (Platform).

On the whole it could be derived that cloud computing does not only pertains to infrastructure, platform or just software as a service but also provide improved and superior capabilities through all the available technological breakthrough till date. The concept of these service provided in the virtualized environment is not at all new, it has been there in the market as parts of grid computing or as web service.

2.3.4. DEPLOYMENT TYPES

Cloud hosting and deployment is totally a different thing to the cloud platform discussed till now. It mostly depends upon the service provider's business model. Previously there has been a propensity of beginning of a cloud to develop from internal solutions as a private cloud and then service is provided to the external users as per their requirements. The private clouds were primarily developed for the fulfillment of domestic requirements, then these internally developed capabilities could be sold publicly for some extra earning. This is how the concept of hybrid cloud took shape in the later part of the cloud evolution. This process of realization of a cloud is thought to be the natural process of growth, but there are no such constraints as such. Thus in conclusion cloud could be defined also by its deployment types narrated below. [49]

- I. **Private Clouds:** These clouds are typically possessed or chartered by individual particular enterprises. The components and functionalities are not transparent to the user except in certain scenarios like SaaS where enhancements are offered according to the user requirements. E.g. eBay.
- II. **Public Clouds.** When enterprises rent or take cloud functionality on loan from other organizations or in case it offers its own personal functionalities to other organizations, the type of cloud is known as Public cloud. This presents the user with the capacity to utilize the different features of the available cloud according to their own requirement or contract out personal requirements reducing the overall cost of

development and ownership. The scope of the cloud functionality varies according to the types of cloud described earlier. E.g. Amazon, Google Apps, Windows Azure.

III. Community Clouds. Basically when small and medium enterprises (SMEs) contribute together to develop certain cloud for commercial purpose and share the profit accordingly. Generally public cloud providers' tender their internal infrastructures as a service to others but in this case normally a single enterprise unable to bear the total input cost of a cloud collaborate with other likeminded organizations to put up a working cloud. These types of clouds are called Community Clouds. Community clouds can either be aggregate public clouds or dedicated resource infrastructures.

This type of cloud could be re-categorized into private or public community cloud. Like few SMEs can pool their resources together to form a private cloud for their own personal use or in contrast somebody like Zimory [145] amalgamates resources from different contributors and resell the services outside.

This type of cloud is still in the process of being fully realized, it has more thing in common with grid computing than in any other type of cloud. Although it still remains a vision in most of the cases but already we do see examples of Community Clouds like Zimory [145] and Right Scale [14].

IV. Hybrid Clouds. In case of public cloud the different enterprises permits the cloud providers a part of their infrastructure, now this causes loss of control over resources, data and even management of code which in most cases is undesirable. Now to resolve such problems the concept of *Hybrid clouds* were conceived. This contains the right mixture of private and public cloud to enable an enterprise to retain preferred degree of control over the data, infrastructure and management as well as outsource the services and attaining capital expenditure.

Today hybrid clouds are very few in the market, the only examples present in the market are IBM and Juniper who have pioneered the base technologies for the realization of hybrid clouds. [4].

2.3.5 . **CLOUD ENVIRONMENT ROLES**

In cloud environments, individual roles can be identified similar to the typical role distribution in Service Oriented Architectures and in particular in (business oriented) Virtual Organizations. As the roles relate strongly to the individual business models it is imperative to have a clear definition of the types of roles involved in order to ensure common understanding.

I. Special Purpose Clouds. In IaaS a cloud, which typically originates from data centers have a “general purpose” influence over the customer. In contrast to this, PaaS clouds tend to provide functionalities which are more specialized to specific use cases, this should not be confused with

“proprietaryness” of the platform: proprietary data entails that organization of data and HMI (Human Machine Interface) is exclusive to the service provider and specialization entails making available extra use case specific technique. [4]

Google App Engine provides specialized functionalities like precise competence for handling of distributed document management. Accordingly it could be anticipated that impending systems would provide further capabilities appealing individual user areas to compete in the ever evolving market scenario.

Special Purpose Cloud is an augmentation over the regular cloud systems which impart extra devoted capabilities which are quite evident in today’s products in the market.

- II. **Cloud Providers** are developers who provide *clouds* to the customer – either via dedicated APIs (PaaS), virtual machines or direct access to the resources (IaaS). It should be noted that hosts of cloud enhanced services (SaaS) are normally referred to as *Service Providers*, though there may be ambiguity between the terms Service Provider and Cloud Provider.
- III. **Cloud Adopters or Software / Services Vendors** are organizations who enhance their own services and capabilities by exploiting cloud platforms from *cloud providers* or *cloud resellers*. This enables them to provide services that scale up to dynamic demands particularly to new business entries who cannot estimate the uptake / demand of their services as yet.

The cloud enhanced services thus effectively become *software as a service*. [137]

IV. Cloud Resellers or Aggregators are people who aggregates cloud platforms from *cloud providers* to either provide a larger resource infrastructure to their customers or to provide enhanced features. These relates to *community clouds* which the cloud aggregators may expose a single interface to a merged cloud infrastructure. They matches the economic benefits of global cloud infrastructures with understanding of local customer needs. They provide highly customized, enhanced offerings to local companies (especially SME's) and world-class applications for important European industry sectors. Similar to the software and consulting industry, the creation of European cloud partner ecosystems provides significant economic opportunities in the application domain by two ways firstly by mapping emerging industry requests into innovative solutions and secondly by utilizing these innovative solutions by European companies in the global marketplace.

V. Cloud Consumers or Users make *direct* use of the cloud capabilities (cf. below) – as opposed to *cloud resellers* and *cloud adopters*, however, not to improve the services and capabilities they offer, but to make use of the direct results, i.e. either to execute complex computations or to host a flexible data set. Note that this involves in particular larger enterprises which outsource their in-house infrastructure to reduce cost and efforts (see also *hybrid clouds*). [137]

It is to be noted that future market developments will most likely enable the user to become provider and consumer at the same time, thus following the “Prosumer” concept, as already introduced by the Service Oriented Architecture concepts [3].

2.3.6. SPECIFIC CHARACTERISTICS / CAPABILITIES OF CLOUDS

Cloud does not present a single individual technology rather it indicates towards a common provisioning model with superior capacity. So it is important to explain in detail these qualities. At the present moment there is a strong inclination towards identifying cloud as a new name to an old idea pertaining to the confusion among the cloud concepts and its relation to the PaaS/IaaS/SaaS models as these aspects have already been attended to without the specific name of cloud being used.

This section specifies a concrete capability associated with clouds that are considered *essential* (required in almost every cloud environment) and is *relevant* (ideally supported, but could be restricted to specific use cases) thereby distinguishing non-functional, economic and technological capabilities addressed, respectively to be addressed by cloud systems.

I. Non-Functional Characteristic

The non-functional characteristics signifies the quality and/or the properties of a system, not the particular technological necessities. The nonfunctional aspects of a cloud could be accomplished by numerous ways and could be deciphered into

many other ways which may create huge issues in compatibility and/or interoperability of different cloud system provider. Today interpretation of cloud differs mainly because of the non-functional aspects. [162,142].

II. Economic Significance

Economic considerations are one of the key primary factors for the establishment of cloud systems in business. The principal advantage of a cloud system is its reduced input cost and effort by means of outsourcing and automation of crucial resources. But the tradeoff lies between the control and the effort. The more loss of control leads to less effort and vice versa. In case of private cloud precise calculations should be made in advance between benefit by cost reduction and the increased effort required for building and running such a system.

III. Technological Aspects

The main technological challenges that can be identified and that are commonly associated with cloud systems are:

- i. Virtualization** is a concept which is much older and may be termed the initiation of clouds which hides the technological complexity from the user and gives enhanced flexibility (through aggregation, routing and translation). More concretely, virtualization supports the following features:
 - a. ***Flexibility and Adaptability***: by creating a virtual execution environment, the underlying infrastructure can change more

flexible according to different conditions and requirements (assigning more resources, etc.).[15]

b. ***Ease of use***: By keeping the complexity of the infrastructure hidden (including management, configuration etc.) virtualization makes it a lot more easier for the user for developing any new applications, also reduces the overhead for controlling the system.

Infrastructure independency: in principle, virtualization allows for higher interaction by making the code platform independent.

c. ***Location independence***: services can be accessed independent of the physical location of the user and the resource.

ii. **Multi-tenancy** is one of the most essential issues in cloud systems, where the location of code and / or data is normally unknown and the same resources may be assigned to multiple users concurrently. This affects infrastructure resources as well as data / applications / services which are hosted on shared resources but need to be made available in multiple isolated instances. Actually all information is maintained in separate databases or tables, but in more complicated cases information may be simultaneously altered, although maintained for isolated tenants. Multi-tenancy implies a lot of potential issues, ranging from data protection to legislator issues[15].

iii. **Security, Privacy and Compliance** is acutely essential in all systems dealing with potentially sensitive data and/or codes.[26]

- iv. **Data Management** is an essential part in particular for storage clouds, in it; data is flexibly distributed among multiple resources. Implicitly, data consistency needs to be maintained over a wide distribution of *replicated* data sources. Same time, the system also needs to be aware of the data location (when replicating across data centers) taking latencies and more particularly workload into consideration. As the size of data could possibly change at any time, data management addresses both horizontal and vertical aspects of scalability. Another important aspect of data management is in providing consistency guarantees. (eventual vs. strong consistency, transactional isolation vs. no isolation, atomic operations over individual data items vs. multiple data times etc.).[26]
- v. **APIs and / or Programming Enhancements:** These are essential for exploiting the cloud features. Common programming models requires the developer to take care of the scalability and autonomic capabilities, while a cloud environment provides the features that allows the user to leave such management to the system.
- vi. **Tools** are generally necessary to support development, adaptation and usage of cloud services.
- vii. **Metering** of all kind of resource and service consumption is essential for offering elastic pricing, charging and billings. It is thus a pre-condition for the elasticity of clouds.[54]

It is up to debate whether the Internet of Things is related to cloud systems at all: whilst the internet of things will certainly have to deal with issues related to

elasticity, reliability and data management etc., there is an implicit assumption that resources in cloud computing are of a type that can host and / or process data – in particular storage and processors that can form a computational unit (a virtual processing platform).

However, specialized clouds may e.g. integrate dedicated sensors to provide enhanced capabilities and the issues related to reliability of data streams etc. are principally independent of the type of data source. Though sensors as yet do not pose essential scalability issues, metering of resources will already require some degree of sensor information integration into the cloud.

Clouds may further offer important support to the internet, for dealing with a flexible amount of data originating due to diverse sensors and “things”. Similarly, cloud concepts in scalability and elasticity might interest for the internet to better cope with dynamically scaling data streams.

After all, the internet may gain a lot of insight from cloud, but there are no direct relationships between these two areas. There are however some similarities that should not be disregarded. Data management and interfaces between sensors and cloud systems displays some common features.

2.4. REVIEW ON DATA INTEGRATION

2.4.1.INTRODUCTION

Data integration is defined as the technique to integrate or collect data from different sources and merge them at one place and finally gives a virtual view to the users. The basic objective of an integration system is to amalgamate numerous information systems into one single unity, so that the user interacts apparently to a single information system. [75]. Integration of information systems seems to be necessary in today's world to meet business and consumers needs. There are two reasons for integration: primarily, to enable a single access point for information exchange and an integrated view. Secondly, towards a particular information need, data from different information systems is accumulated to gain a more comprehensive basis towards the required need.

There are so many applications that are benefited from integration. In Business Intelligence integrated information is required for generating different reports and for executing queries. The CRM (Customer Relationship Management) is another example of a field where extensive requirement of integrated information is required to evaluate individual customers, business trends, sales, etc. to improve customer satisfaction.

Integrated information of a company is presented in a personalized website which are normally called Enterprise Information Portals (EIP) which provide a single point of interaction for any information interchange. Lastly, in the area of E-

Commerce and E-Business, an integrated information system acts as a facilitator as well as an enabler towards business transactions and services over computer networks [63].

Integrating multiple information systems creates a unified virtual view to the user's imperative of the number of system or location of the actual stored data. The users are offered a standardized logical view, which in reality is scattered over diverse data sources. To enable such a scenario the same data abstraction standards has to be maintained (unified global data model and unified semantics)[89].

Data integration is hard. The evidence is overwhelming. Every company we've talked to about their data has data integration problem. It's not just the IT people that moan about it either, it's IT users too and the company executives. Everywhere data is almost in a constant mess throughout. Today we have a dedicated sector of the industry devoted towards data integration solution; it generates about \$3 billion in revenue and its growing space. Aside from that there are probably billions more spent on in house data integration efforts whether they employ the whiz data integration tools or not [89].

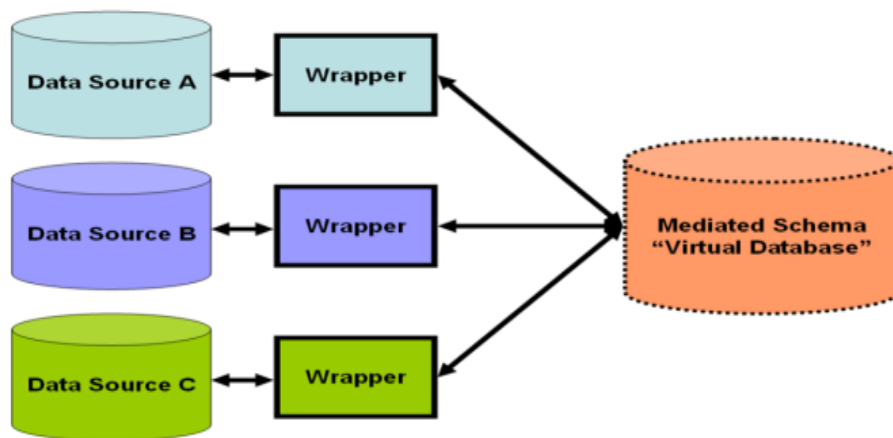


Figure 2.7: Data integration chart

2.4.2. CHALLENGES OF DATA INTEGRATION.

Basically, an information systems design normally does not keep provision for any type of integration, so the main job is to identify data and schema conflicts while integrating. So for integration adaptation and reconciliation functionality has to be applied to unite the conflicting data and schema. While the goal is always to provide a homogeneous, unified view on data from different sources, the particular integration task may depend on the following [161]:

I. Architectural View:

It is the art of expressing a model or a concept of information for utilization in activities which requires explicit details of complex systems. The most common activities in this sector are library systems, Content Management Systems, Web development, User interactions, Database development, Programming, Technical writing, Enterprise architecture and critical systems software design. As different information systems have different architecture, it becomes very difficult to integrate information from different architecture. [161, 75].

i. Content and Functionality of the Component Systems:

Component systems of any architecture have contains different content and have different functionality which made it difficult to integrate [161].

ii. Source Wrapping:

To access data sources on the Web, a crucial step is wrapping, which translates query responses, rendered in textual HTML, back into their relational form. Traditionally, this problem has been addressed with

syntax-based approaches for a single source. However, as online databases multiply, we often need to wrap multiple sources, in particular for domain-based integration.

iii. **Semantic Data:**

A chief requirement of data integration systems is that the differences in the syntax and semantics of the underlying data sources should be hidden from the user. However, data sources are often independently created and do not easily interoperate with other sources.

The rise of XML as a data interchange standard has led to a new era in which research has focused on semi-structured data. Although XML has provided a single interchange format, different users can model the same data in different ways, which can lead to heterogeneities at various levels, including the semantic level. Semantic heterogeneities can arise from entities being perceived differently.

iv. **Streaming Data:**

Now a day in scientific and industrial environments, the amount of data in form of heterogeneous streams is becoming one of the main sources of information and knowledge acquisition. Advances in wireless communications and sensor technologies have enabled the deployment of networks of interconnected devices capable of ubiquitous data capture, processing and delivery of such streams [89].

v. **Large Scale Automatic Schema Matching:**

In the integration of heterogeneous data sources, schema matching becomes a critical problem. Basically as well as traditionally, the problem of matching multiple schemas has essentially relied on finding pair wise - attribute correspondences in isolation. Schema matching is a basic operation of data integration and several tools for automating it have been proposed and evaluated in the database community. Research in this area reveals that there is no single schema matcher that is guaranteed to succeed in finding a good mapping for all possible domains, and thus an ensemble of schema matchers should be considered. Informally, schema meta-matching stands for computing a “consensus” ranking of alternative mappings between two schemata, given the “individual” graded rankings provided by several schema matchers [114].

vi. **Construction of Global Schema:**

Schema here referred to structure data residing at a place has a structure different than other and integration of data from different structure proves to be difficult.

- a. Kind of info-this include alphanumeric data, multimedia data; structured, semi-structured, unstructured data.
- b. Available resources- time, money, human resources, know-how, etc.

vii. Understand Data Needs:

It could be defined as the delivery of the right data to the right application in order to achieve the right business result. Primarily this is the main reason for which we tend to formulate corporate data centers and to ensure data moves to the appropriate location. Every change in data structure of any single unit has impact on the whole architecture of it.

viii. Understand Business Timing Needs:

The major contributing factor to any data process is the actual business activity for which it is required. Data is the most important asset owned by any company. It is IT's almost sacred duty to deliver the data where it is needed when it is needed [89].

Data may have many target systems, some need the data in real time others, like the BI system, only require periodic updates. Therefore the integration solution must be able to handle batch and real-time activity.

I. Integrate Master Data and Governance Rules:

Where MDM solutions have been implemented then the MDM becomes the Centre of the hub for particular types of data. e.g., all customer data must be validated against the customer master. In this case customer data must be validated by the customer master before being forwarded to other systems that require the data. Thus data distribution can be a two-step process.

II. Technical:

Some of the most significant technical challenges of designing an application integration environment involve identifying the technical needs of your solution and determining the combination of protocols and services that will provide for those needs.

III. Organization Issues:

An application integration environment is available in multiple departments in an organization. Staff in different department may choose to deploy application that will need to integrate with your application integration environment.

There may be some kind of heterogeneity which includes differences in-

1. Hardware and operating systems
2. Data management software
3. Data models, data semantics and middleware

2.4.3. APPROACHES TO INTEGRATION

I. Integration by Application:

Process of bringing data from multiple application programs into an unified program. It is secured and orchestrated approach towards transferring & sharing of processes or data between different applications within the enterprises. This is applicable for small number of component systems and application become as the number of systems interfaces and data formats integrate [161].

a. Benefits:

Application integration allows the application to be introduced into the organization more efficiently for faster work and more accessibility at a lower cost. It allows you to modify the business processes as per the requirement of the organization. Providing so many channels for organization where they interact with each other and get integrated.

b. Technique:

There are two models which increase the efficiency of application integration.

i. Point To Point Model:

It is decentralized structures in applications communicate directly with each other applications. This is most useful for organizations where they have few applications with small number of sources.

ii. Integration Hub Model:

In this an integration hub is placed between the applications at there each application is interact or communicate with it rather than communicating with each other. Application needs only a interface and connection to the integration hub, when a new app are introduced you do not need to rewrite the interface.

c. Requirement for Application Integration:

It provides common interface through which application can communicate with each other. There must be strong connectivity between the platforms to avoid any

disturbances. A common set of process and services rules should be used to insure consistency and reuse of integration services. Be capable of reusing the existing transport protocols that already exist in the enterprises.

II. Common Data Storage:

Data integration can be done by transferring data into new data storage. Usually it provide fast data access and easily understandable by users. If in some case local data sources are retired or damage, the application used by them is moved to new data storage .a common data storage have been refreshed periodically so that the local data sources remain functional and easily accessible for the user [75]. This type of integration is successful in organization as it increases accessibility among the people working in any organization.

a. Benefits:

As data is stored is stored at one place so it is easy for data mining. It reduces the cost and time to produce input formats provide the basis for consistent database reuse. It reduces the risk of errors caused by formats conversion failures. Data storage makes it easier to test and compare results of different federations.

b. Technique:

Data warehousing is an approach towards realization of a common data storage and integration. Data is extracted, transformed, and loaded (ETL) into a data warehouse from several mostly heterogeneous operational sources. On this data different analysis tools could be implemented e.g. OLAP, OLTP.

Another common storage example is Operational Data Storage. In it “Warehouses with fresh data” is constructed by “Real Time” updates in local data sources to the data store. This further enables a real time data analysis process for decision support systems. Opposed to data warehouses, data here is neither cleansed nor aggregated.

III. Uniform Data Access:

It is defined as connectivity and controllability across various data sources in this data from all the sources which are developed in different structures, schemas, and architecture are accumulated at one place which is treated as virtual data. Since it is a time consuming process, thus data access, homogenization and integration have to be done at the runtime. Without it user face the difficulty to translate various data sources into the format supported by their application to provide a collective single view of data.

a. Benefits:

It is ideal for application developers, software vendors and solution providers. Increase in the consistency in data accessing, accessing of become easier anywhere in the enterprise. Open to each and every person within an organization. Editing can be done to a data without interrupting the normal state of data. It can improve the use of all assets – hardware, software and people.

b. Techniques:

In this process mediated query system is implemented as a solution to a single point for read only querying access to various data sources. These sends sub-queries to local data sources and are then combined at local query results.

P2P integration (peer to peer) –It is a decentralized approach where data is mutually shared and integrated at every peer location. This is entirely dependent on integration functionality available at all the peer location.

FDBMS- Federated database systems (FDBMS) try to achieve a uniform data access solution by integrating data from underlying local DBMS into a logical layer. Federated database systems are fully functional DBMS. They normally have their own data model, global transactions, global queries support as well as global access control. Usually, the five-level reference architecture is employed for building FDBM [77].

Portals are personalized access ways to an uniform data access for information on the internet or intranet where each user is provided with information tailored to his information needs. Usually, web mining is applied.

IV. Manual integration:

Here users are directly interacting with relevant information systems and integration can be done by selection of data from various data sources. As data has been developed at different structure and architecture it seems to be easy. We have to deal with different user interfaces and query languages. Detailed

knowledge on location of data is necessary as we have to take specific data for integration. A logical data representation needs to be very accurate while doing it. Logical data length can be 1, 2, 4 or 8 bytes in length. One another important thing Data semantics should be there means connection of database to the real world outside the database or mapping between an object modeled represented and/or stored in an information system.

a. Benefits:

The basic benefit of this type of approach is its accuracy, and adjustability to any type of requirement. In any place where high resolution and accuracy of data is required or the data has an irregular pattern it is the most convenient and appropriate approach for data integration.

b. Techniques:

It requires writing of large numbers of programs/queries, or some existing programs are required to be modified for further customization of the program as the requirement may be. The most labor intensive part of manual integration is data mapping. In data mapping a point to point approach is required, which an operator has to manually find and map different data, before the integration could be possible. All this can lead to increased complexity with the addition of every new data source. On the downside it costly, time consuming, and error prone.

Since the inception of the concept of DBMS one of the major problem which has been dogging database design is to build a database structure with respect to the

heterogeneous user environment, applications and diverse data requirements are satisfied simultaneously.

For manipulation and operation of database the management system requires a definition, this definition is stored in the form “Schema” This schema is also defined as the intension of the database. Instances or occurrences are the actual values of data in a given database.

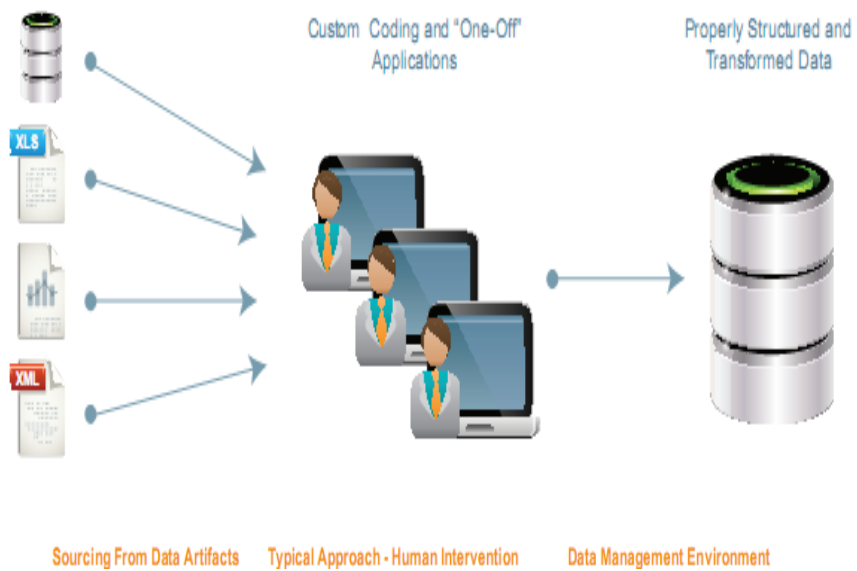


Figure 2.8: Manual Data Integration Process

2.5. INTRODUCTION TO EXPLORATION & PRODUCTION DATA

ORGANIZATION.

The advancement in information technology has provided with an exceptional volume of data. As a matter of fact data has become the life sustaining paradigm of any oil and gas company today. From allocation of sales, production to designation of the location of the any drilling all the decisions are established on pure data. But in case of realization of competitive advantage, the concern is not the amount of data available to a company, but how efficiently it is able to exploit the available data to gain business intelligence and operational knowledge. Thus data management is a paramount concern extending over every discipline within the modern exploration and production establishment.

2.5.1. BASIC DATA INTEGRATION CHALLENGES IN E&P SECTOR.

With few exceptions, respondents agreed on a common set of challenges facing integration solutions, including:

- I. Absence of naming conventions and universal standards;
- II. Information Silos;
- III. Absence of a common integration framework
- IV. Inconsistency in quality and gaps in data;
- V. Enormous amount of data;

The above mentioned problems of inefficiency, obstacles not only have technical but also cultural aspects. The prevalent problem which dogged the entire E&P

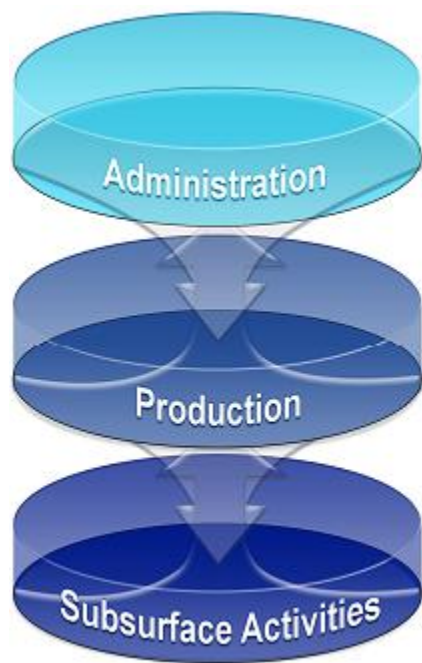
industry since its inception is the amount of data and silos of those data. Another major problem which is contributing to proper integration of data is the standardization of naming convention. In this research the main aim is to standardize the main data integration architecture so that seamless data integration occurs imperative of their location, heterogeneity or place of origin.

Upstream segment of oil & gas industry are basically divided into three primary areas that mutual collective information interest and consequently a requirement to exchange information. Those three domains included:

i. Administration: Consisting of Accounting, management reporting, regulatory compliance – Utilizing production information to realize tactical and long-term strategic business decisions.

ii. Production: Consisting of Production, accounting/allocation, Process control, and Asset management – Applying real- time information to enhance production, better manage assets, and arrange well maintenance.

iii. Subsurface activities: Consisting of Seismic, well analysis, modeling – Applied on geophysical and geological (G&G) information.



2.9. Upstream Domain Classification

There are sources of value in exposing production information to stakeholders across these three functions, for a variety of good reasons.

The Exploration and Production (E&P) sector can be divided into three basic decision making scenarios. It normally depends upon the response time taken to act after a certain incident. The incidents which might occur are divided into two categories namely (i) Planned and (ii) Unplanned in every decision scenario. The planned occurrence is well thought out procedure which has a greater time to respond and thoroughly planned to a schedule. Unplanned occurrence of any event is much more complicated and could lead to huge monetary losses or can even lead to loss of life. The more is the response time the amount of loss is increases proportionally, or if a wrong decision is taken in a hurry it might aggravate the situation. So it is highly prudent to have a system which would be able to provide correct information to the correct person at an optimized time.[70]

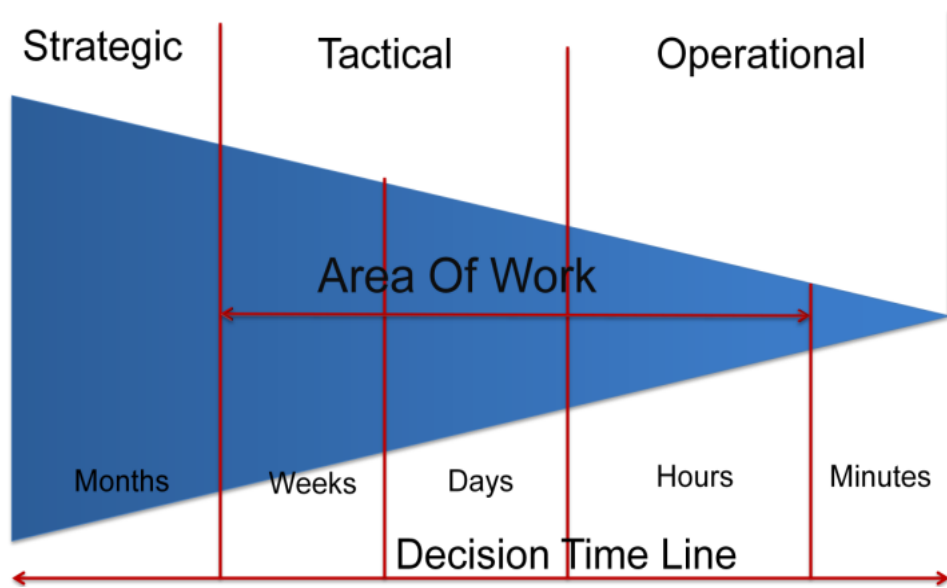


Figure 2.10. Area of work

Namely they are defined as operational, tactical and strategic.

	Data Type	Volume of Data
Operational	Highly Structured	Detailed and Huge
Tactical	Structured & Semi-structured	Medium
Strategic	Semi & Unstructured	Highly condensed

Table 3.1: Operational Data Types[70]

Operational:

Data generated and processed during the day to day operations of exploration and production centers. These operational decisions have a really short response time ranging from minutes to hours. The table given below describes the operations which occur in the process.

	Operational Activity	Activities Leading To Data Handling
Exploration	Remote sensing survey	Remote Sensing Satellite & Low Flying Aircraft Images
	Seismic survey	<ul style="list-style-type: none"> ➤ Onshore sites and marine resource areas ➤ Possible onshore extension of marine seismic lines ➤ Onshore navigational beacons ➤ Onshore seismic lines ➤ Seismic operation camps
	Exploratory Drilling	<ul style="list-style-type: none"> ➤ Access for drilling unit and supply units ➤ Storage facilities ➤ Waste disposal facilities ➤ Testing capabilities
Production		<ul style="list-style-type: none"> ➤ Wellheads ➤ Advanced recovery techniques ➤ Flow lines ➤ depleted Separation / treatment facilities ➤ Oil storage ➤ Facilities to export product ➤ Flares ➤ Gas production plant

Table 3.2: Operations for Exploration & Production [70]

Tactical.

Tactical decisions have a response time varying from days to weeks. It mostly consists of minor to semi major problems of day to day activities, but impact of those decision needs authentication from local administration to central administration.

Since the inception of the concept of DBMS one of the major problems which have been dogging database design is to build a database structure with respect to the heterogeneous user environment, applications and diverse data requirements are satisfied simultaneously.

For manipulation and operation of database the management system requires a definition, this definition is stored in the form “Schema” This schema is also defined as the intension of the database. Instances or occurrences are the actual values of data in a given database.

Logical schema contains logical specifications like groupings of attributes and relationships among these groupings. Whereas as physical specifications like ordering, access to records, indexes, and physical placement. The database design is distinguished based on these physical and logical schema design. Logical schema design involves the problem of designing the conceptual schema and mapping such a schema into the schema definition language of a specific DBMS. In Figure 1 the phases of database design is represented in graphical format. The different phases of database design are as follows.

1. Requirements Specification and Analysis.

The primary specifications of the requirements of various information for a particular person, user groups or an organization is analyzed for identification and quantification are done.

2. Conceptual Design.

Conceptual design of a schema is modeled in this phase to represent the applications' and users' views of the required specification of the processing or utilization of information.

3. Implementation Design.

Conversion of conceptual to logical schema of a DBMS is described as implementation design. Logical database design is described as the result of second and third phases.

4. Physical Schema Design and Optimization.

For optimization of performance in accordance to a set of transactions the logical schema of a database is mapped to stored representation, this is call physical schema design and optimization.

Typically, the application design activity proceeds in parallel with database design. Hence, Figure 3.3. also shows specifications related to applications as the outputs of the last two phases. As shown in Figure 3.3, the activity of view integration can be performed at several points of the database design process. It usually is performed during conceptual design. In that case, its goal is to produce an integrated schema starting from several application views that have been produced independently.[70]

2.5.2. TYPES OF DATA INVOLVED IN E&P SECTOR.

I. Geological Data:

Geology is basically the study of earth's features and properties. In Greek "gê" means Earth and "logas" means study, it is the study of rocks for its composition their process of change of properties e.g. igneous rocks after a long time of exposure to pressure and temperature changes into metamorphic rock which has totally different properties than the former igneous rocks.[161]

Geological study is paramount in case of finding hydrocarbon, mineral or for detection of ground water. Geological data may consist of bottom sampling, shallow coring, stratigraphic analysis, rock formations etc.

Rock Cycle

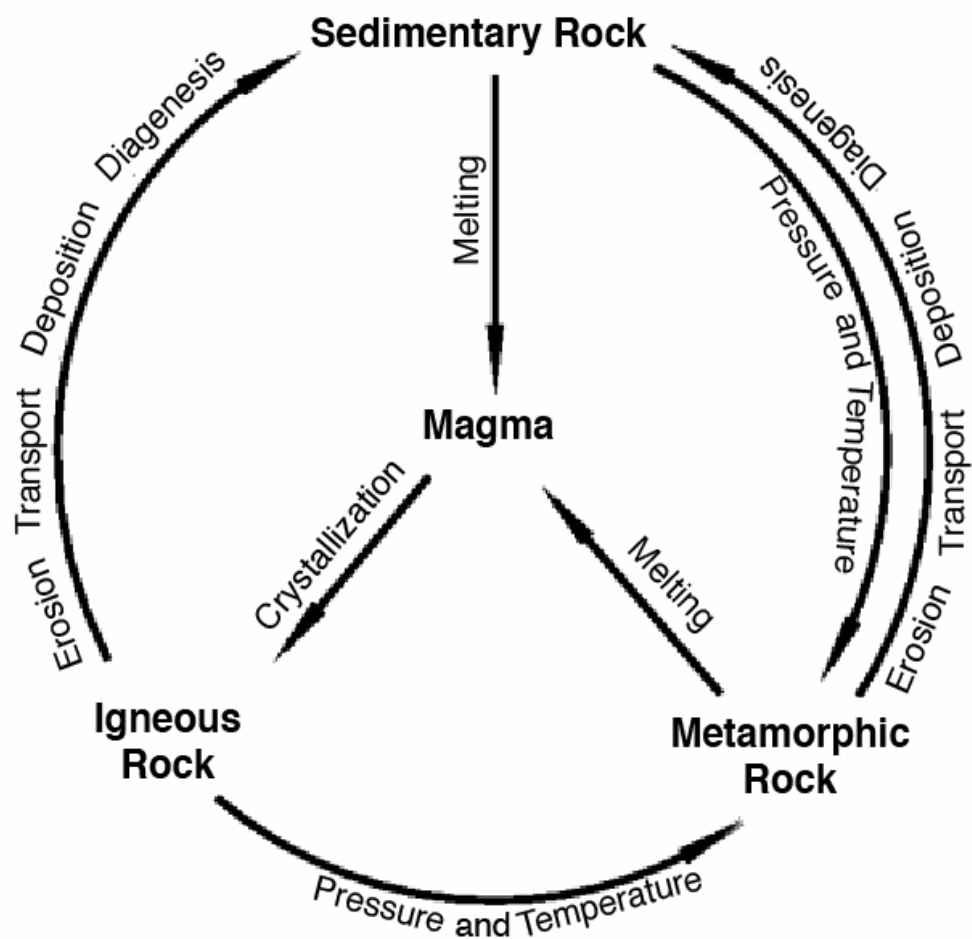


Figure 2.11. Rock Cycle [161]

II. Seismic Data/Geophysical data:

Seismic data is acquired through a process of seismic survey. This is based on the physics principal of reflection of sound wave to determine different layers of earth's subsurface layers. In this method a controlled artificial seismic wave is generated through explosion of dynamite, "Tovex" a type of air gun or a thumper machine which converts mechanical energy into seismic wave called "Vibroseis". [151]

Seismic wave are basically synthetic mechanically generated sound waves which are regulated by acoustic impedance of the medium through which it travels, in this case which is Earth. The acoustic (or seismic) impedance, Z , is defined by the equation:

$$Z = V\rho,$$

where V is the seismic wave velocity and ρ (Greek *rho*) is the density of the rock.

The earth crust is made of different layers of rocks each of which has a different density, when acoustic or sound wave is incident on these layers some part of the wave penetrates the layer and the other part gets reflected. These reflected rays are received by devices known as "Geo-phones" placed at strategic locations. In this process the components like time taken by the waves to travel from source to destination is measured, the velocity of the acoustic wave is known and with these knowledge the image of the subsurface is reconstructed by the geophysicists.

Other types of geophysical data consist of gravity and magnetic survey. These are conducted by gravity meter and magnetometer respectively. The magnetometers are of two types, with a single sensor for the measurement of the total magnetic field strength and with multiple sensors for gradient measurement of magnetic field.

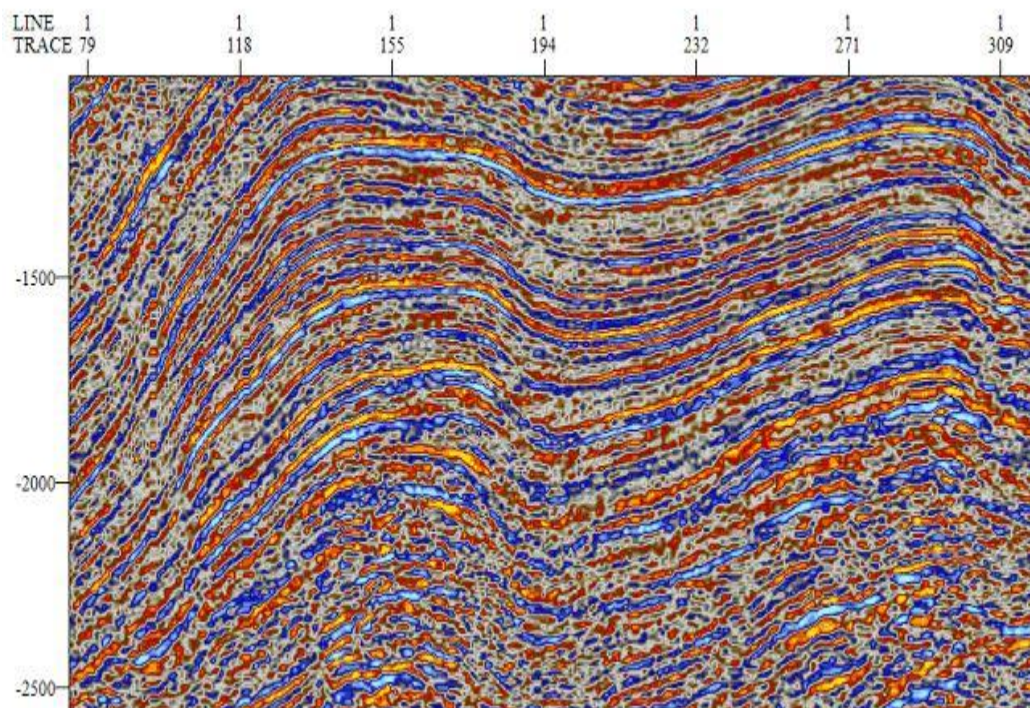


Figure 2.12. Seismic Survey Image [151]

III. Drilling Data

Drilling is a precise process by which an insertion is made into the earth surface to find hydrocarbon. When a potential geological structure is recognized, confirmation of the presence of hydrocarbon is possible only through drilling a borehole into the earth surface, this type of bore hole is known as “exploration well”. In drilling process the different types of data which are required to be collected and processed in this method of drilling consists of Directional drilling / Geosteering, Drilling fluid invasion, Tracers, Underbalanced drilling, etc. [161,151]

IV. Well Log Data

The comprehensive documentation of the geological formations infiltrated by a bore hole is known as “Well logging” or borehole logging. The record or log could be generated by two different processes, viz. Geological (when the log is generated through visual assessment) and Geophysical (when log is generated through measurement of physical properties of the sample by instruments). Well logging can be done during any phase of a well's history; drilling, completing, producing and abandoning. Well logging is performed in boreholes drilled for the oil and gas, groundwater, mineral and geothermal exploration, as well as part of environmental and geotechnical studies.[152]

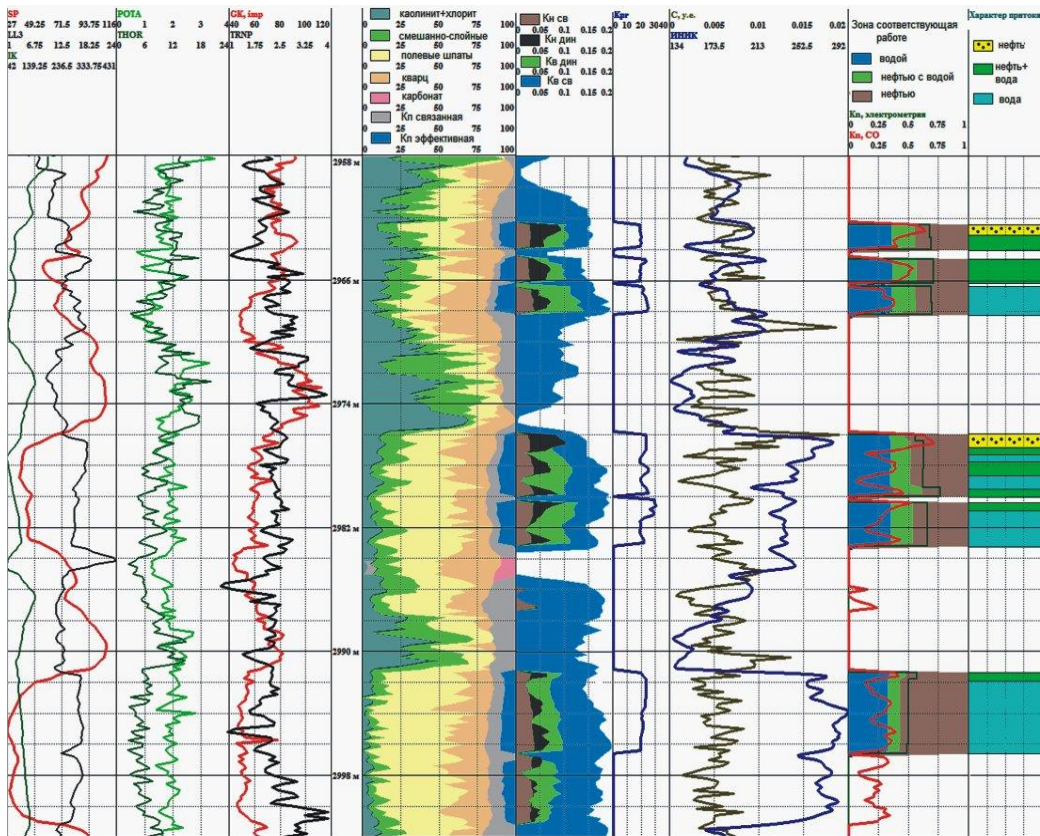


Figure 2.13. Well Logging Data [152]

The different types of logs maintained for this process are as below:

- a) Gamma ray logging
- b) Spontaneous potential logging
- c) Resistivity logging
- d) Density logging
- e) Sonic logging
- f) Caliper logging
- g) Mud logging
- h) LWD/MWD (Logging while Drilling/Measurements While Drilling)
- i) NMR Logging (Nuclear magnetic resonance)

V. Production/Reservoir Data.

Naturally occurring hydrocarbons sometimes accumulates inside porous rocks inside the subsurface of the earth forming a pool, this is addressed as petroleum reservoir, or oil and gas reservoir. Reservoirs are found using hydrocarbon exploration methods. [152]

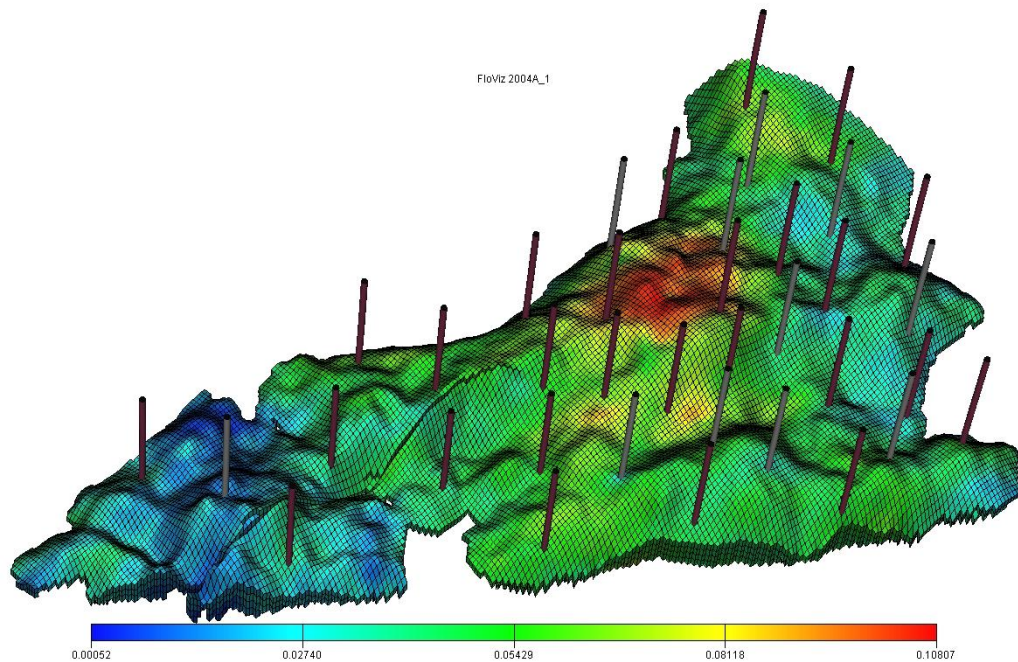


Figure 2.14. Reservoir Simulation Image [152]

2.6. REAL-TIME DATA INTEGRATION.

2.6.1. INTRODUCTION

“Real-Time” doesn’t have a definite or a single definition so it entirely depends on the realization of the requirement of the business an organization is handling.

‘Real time’ could be defined in the following ways:

1. To achieve zero latency in a process
2. Any time access to information by a process whenever required.
3. The process is able to provide information to the management whenever and wherever it is required.
4. To be able to derive key performance measures which communicates to the situation at the present point in time and not just to some historic condition.

Taking into consideration of the above given point a real time information system enabling traditional decision support but the extraction and processing of data from the operational management systems is done with zero latency and Based on these descriptions, RTBI provides the same functionality as the traditional business intelligence, but operates on data that is extracted from operational data sources with zero latency, and also makes it possible to provides means to promulgate measures back into the business processes in real time.

Explicitly, real time information system comprises of:

1. Real-Time delivery of information,
2. Real-Time modeling of data,
3. Real-Time analysis of data,
4. Real-Time implementation of decisions derived from the available information.

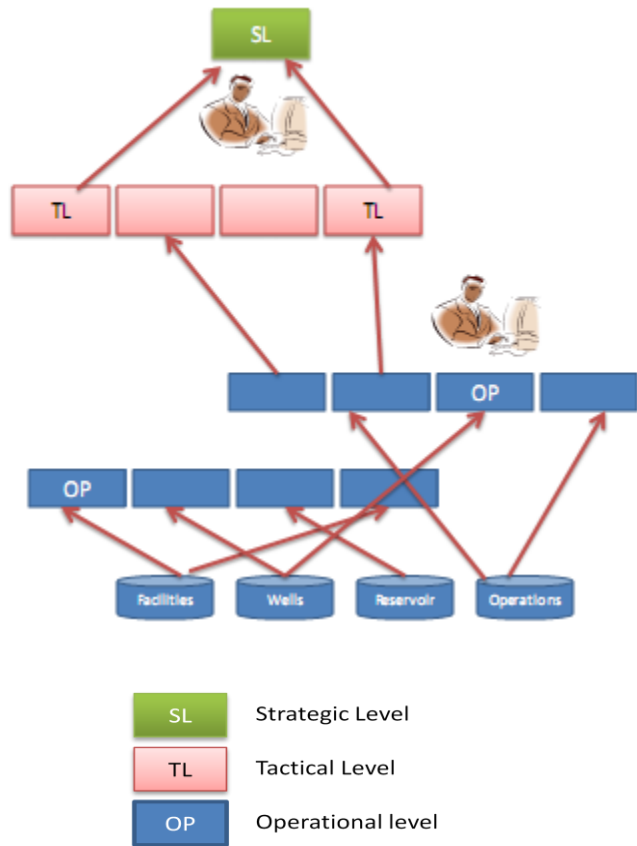


Figure 2.15. Current organization of data

2.6.2. CURRENT WORKFLOWS AND PROCESSES INTEGRATION

Almost every chief E&P Company follows analogous workflow process for transmitting production data en route for office. Additionally there are homogeneous array of infrastructure and applications which support are aligned with projected workflows to acquire and supply production data into enterprise systems like SAP or JD Edwards. The main objective is to augment automation and diminish time and labor costs from work processes while obtaining timely production information from the field.

.An alternative way of classification is bifurcation of technical and financial sectors. Technical systems consist of Real-time systems like SCADA and DCS which pertains basically to production operations. This information is utilized to plan and schedule well workovers. The financial department hardly requires this technical information but needs to concentrate on the market scenario and price dependency derivatives.

In the next chapter the designed architecture and algorithm is given as a solution of the problem discussed in the first two chapters of the thesis. It also contains the testing scenario and discusses the results of the software developed for testing purpose of the developed algorithm.

CHAPTER 3

REAL TIME DATA INTEGRATION ARCHITECTURE

(RTDIA)

This chapter defines the design and components of the proposed architecture for the fulfillment of requirements to attain a real time data integration environment. The later part of the chapter deals with the algorithm developed for the identification and integration of data from remote data bases into the central data warehouse and the software developed for testing of the algorithm.

3.1. REAL TIME DATA INTEGRATION ARCHITECTURE.

The basic idea behind designing a “Real Time Integration Architecture” is to provide flawless conversion of data to information and that back into action. In the current E&P data integration and processing sector two major bottle-necks have a slowing down effect on the system which creates a barrier for achieving the goal. Contemporary architecture suffers primarily from improper data identification and then an inefficient data integration system in the E&P sector which leads to a huge loss of time, skilled manpower and sometimes to loss of precious natural resources. Although the E&P industry has highly skilled human resource available for the analysis part but the above stated problems of data integration are a constant source of hindrance.

Figure 3. exemplifies the current scenario of the E&P industry's information management system. RTDIA's basic objective is to diminish resistance and enable data integration between operational, tactical and strategic levels of process control. Fundamentally the two objective of the architecture would be:

I. Capability to deliver value addition.

In the current scenario companies are not only hindered by less resource, extremely constrained time line but also has to deliver higher throughputs. To dispense enhanced outcome faster, G&G and engineering departments requires deliberating on their domain works and not on finding and processing of data. Highly skilled workforces are costly to come by, so it is vital that their work hour is properly utilized.

II. View Integration.

To understand the proper and correct working scenario requirement of an integrated view which divulges pertinent information at a single view is overriding. e.g. in an exploration situation, the comprehensive outlook must comprise of tornado analog well production histories, charts that measure risk, rock properties, rig schedules, log files, and other variables relating to the prospect in question.

The concept of this research is to implement recent technological enhancement to device an architecture which would be able to eliminate manual intervention existing in current systems and automate both the flow of information from

operational to tactical to strategic layer, representing data to the information stage, and the actions necessary to translate strategic objectives back to operational drivers to effect strategic decisions in real time.

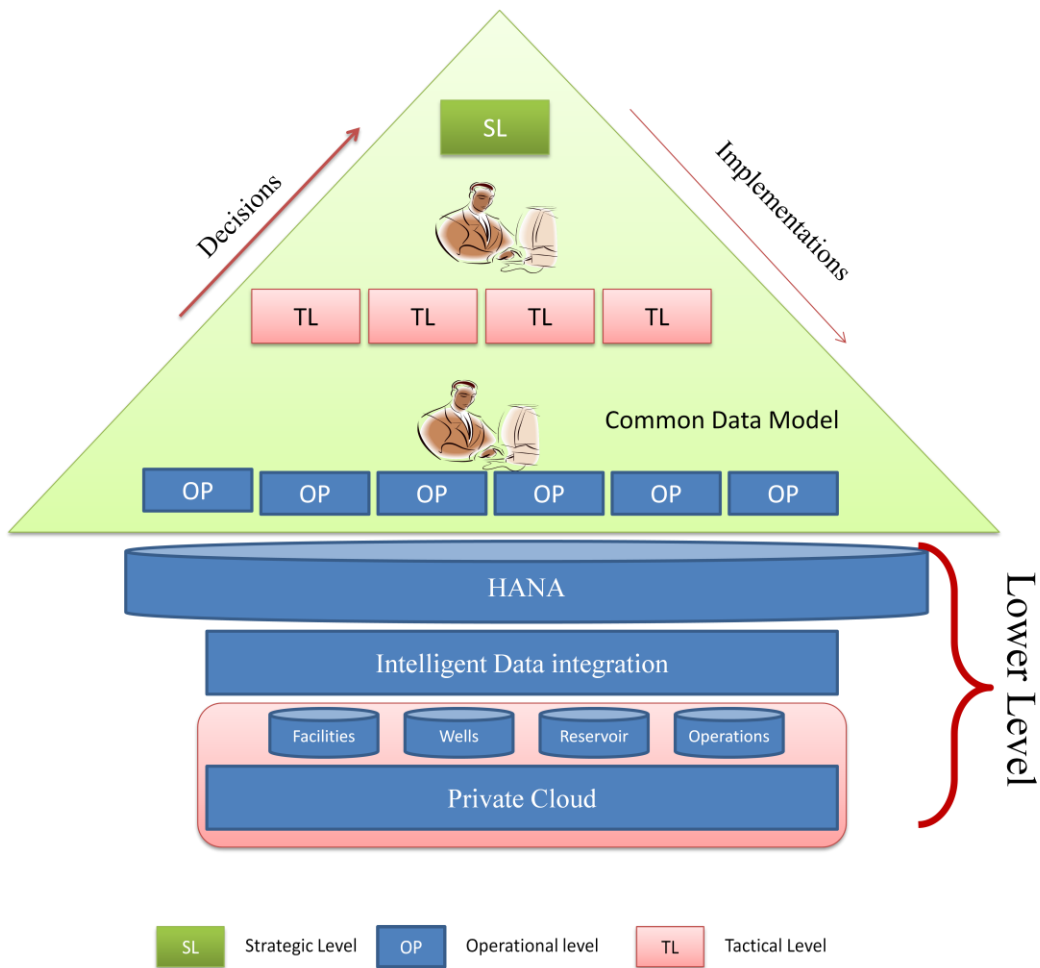


Figure 3.1. Architecture I [Upper Layer].

The expanded form of the lower level of the architecture is described in an expanded form in figure 3.6. in this chapter.

The above given figure describes a real time data integration architectural overview the more prevalent integration tools are basically inert in nature and their sole purpose is to accommodate reports and dashboard. Current BI tools are mainly passive and their objective is to cater for the information consumer by providing reports or at most dashboard-like monitoring of various business processes. With an real time architecture in the information management system which allows effective and seamless integration throughout heterogeneous systems the process of data organization and extraction of operational information for decision support systems would become a effortless procedure and proper attention could be given to functioning of control processes. The architecture needs a smart algorithm to converge the data flow from the physically distributed system. The algorithm needs to be able to efficiently identify the requirement from intelligent metadata repository as well as a rollback schema as a failsafe backup.

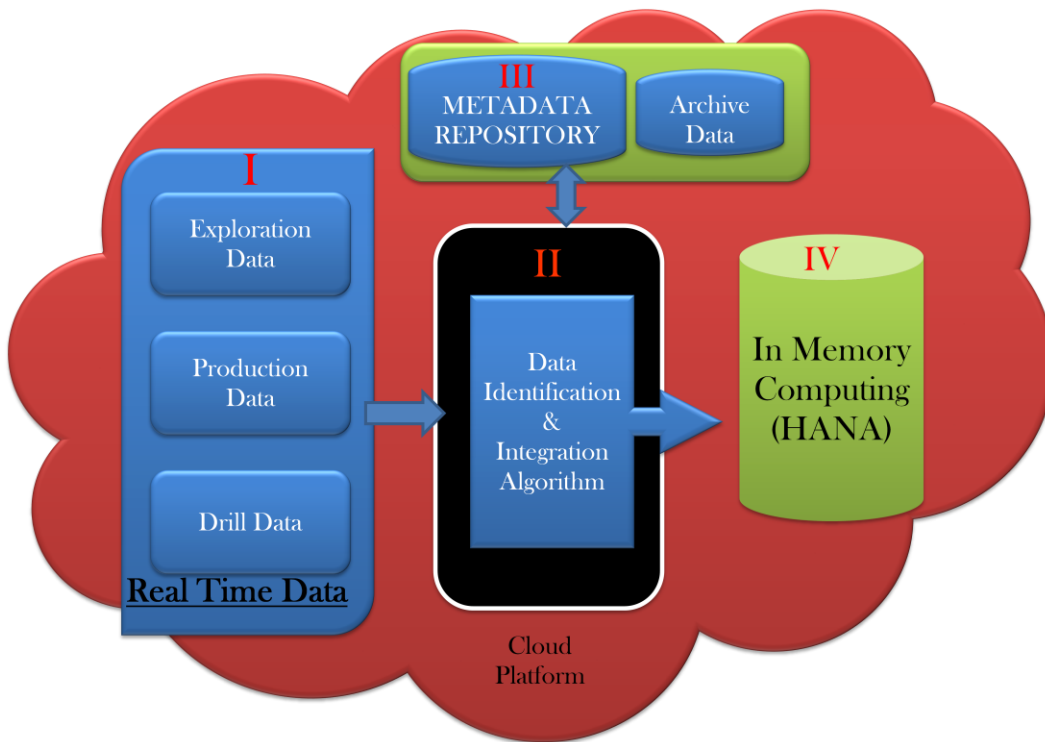


Figure 3.2. Expansion of “Lower Layer” Architecture

3.2. THE LOWER LEVEL ARCHITECTURE (EXPANDED).

The best solution for any problem is the simplest one. In the above depicted architecture the base platform of all real time data generation like SCADA system is based over a cloud platform, which in turn will provide a highly virtualized environment to enable unimpeded communication among the different heterogeneous data sources. With cloud storage platform the most inconvenient problem which has dogged illustrate E&P industry till date, of heterogeneous source databases for integration would invariably reduce the time and effort for preprocessing of data required before diverse database integration.

There are two types of logical data marts designed in the architecture above

I. Dynamic Logical Data Marts. (*Part-I of the figure*):

As the name illustrates, the data marts/ specific sector data warehouses which varies with time is defined here as dynamic logical data marts. The time variant databases are the real time data generations processes where the type and definitions of data remains constant over time but the value changes with time which is ultimately required for any further processing for Business Intelligence. These databases contain just the day to day values generated by the different operational processes which transpire everyday in course of exploration and production of hydrocarbon. e.g. exploration data marts, production data marts, drill data marts.

II. **Static Logical Data Marts.** (*Part-III of the figure*)

The specific databases which remains constant over time like the definition/schema of a database is defined as static logical data marts.

They basically consists of two parts the

- a. Meta Data Repository
- b. Archival Data Repository.

Although these databases are not static over time in the real sense but the time span before which anything changes in these databases occurs is quite big, so these databases can safely be affirmed as static.

a. Meta Data Repository:

It contains critical information about the data units which are stored in the dynamic data marts which works in real-time. The basic definition of every data units has to be well defined without any ambiguity as this could lead into data identification problem which will further cause significant loss in time for data integration which is the main aim of the research.

b. Archival Data Repository

It is separate database created for storage of historical data for reference purpose, or for different trend analysis. The storage could be in different media thus the identification and extraction of these data could vary from seconds to hours

depending upon the category of information as well as the type of media it is stored in. The storage media normally could be detachable magnetic cartridge, online hard drive or robotic tape drives. The data is segregated according to the requirement and importance and divided into these different storage media.

3.2.1. DATA INTEGRATION LAYER (*Part-II of the figure*).

The main conception of this architecture is the smart data integration layer. Customarily this amount of data is handled through ETL processes which are an extraction, transformation and loading tool which wheedle out data from different data marts and load into the central data warehouse for any type of analysis purpose. But ETL tools have their own pros and cons as discussed in chapter 2 of this research report. The designed algorithm as stated later in this chapter works in simple steps by identifying the data vertical by its asset-id, then brings in all the related data to the HANA database and creates a table in the In-Memory. From the In-Memory database table different slicing and dicing of data is performed according to the requirement

I. Benefits:

1. Unified Data Layer:

A collective meta-data repository configuration amalgamates data access by establishing a virtual warehouse view of every required distributed data so

that users, in spite of their departments, have entrée to the same values and sources.

2. Streamlined Processing Cycle:

Logical data integration through a source description generator brings a simple streamline structure to the flow of data from the heterogeneous data sources to the data warehouse. The link created between a unified meta data repository and the data source makes it absolutely seamless processing and integration system.

In this designing attempt effort was to design an optimum architecture which would be able to handle the humongous pressure created in an E&P company and still be able to provide data in a real time scenario. The heart of this architecture is the data identification and integration algorithm.

3.2.2. HANA DATABASE.

Semiconductor memory is by far the fastest type of memory available today and with the increasingly depreciating price of semiconductor main memory (RAM) it is now financially feasible to have large arrays of RAM for high speed processing of colossal volumes of data. HANA Database enables the user to utilize the high speed of semiconductor drives as well as multi core processor to give ultimate proficiency in data analysis and transaction. The SAP HANA supports relational data, graph and text including semi and unstructured data as well contained in the same system making it the most effective option for the current problem. SAP

HANA is also cloud enabled which makes it easier to incorporate in the above described architecture. HANA being 100% ACID(Atomicity, Consistency, Isolation, Durability) compliant making it one of the most reliable Database. [46,116]

In the below given figure 3.3 and figure 3.4 the HANA architecture is depicted in two stages. The figure 3.3 describes the base architecture of HANA database. In its basic form, it shows three basic components of In-memory computing. The first part consists of the distributed databases from where the data is retrieved and brought into HANA. The third component is the BA/BI tools which communicate to the high speed memory to for further processing of data.

Figure 3.4 depicts the different components which and how it communicates with the HANA engine. Basically the distributed databases are accessed by the developed real-time algorithm as well as the metadata repository through high definition file servers. The extracted data is then brought onto the HANA engine through a Sybase replication engine without making any changes in the primary data source and creates a table which is accessed by any integrated BI tools, making it highly agile and effective tool for real-time data operations.

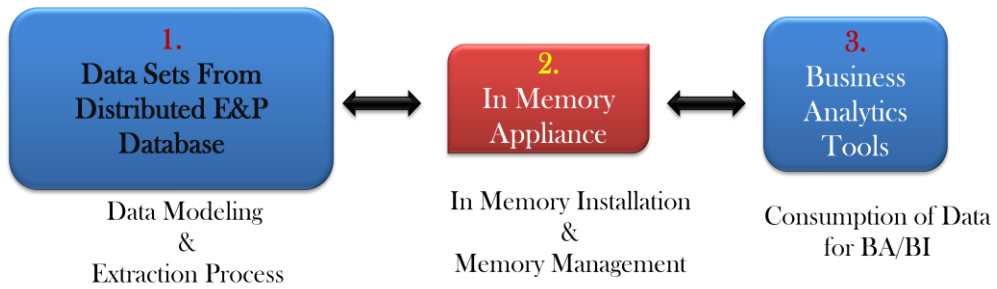


Figure 3.3. Basic Data Flow Diagram of HANA

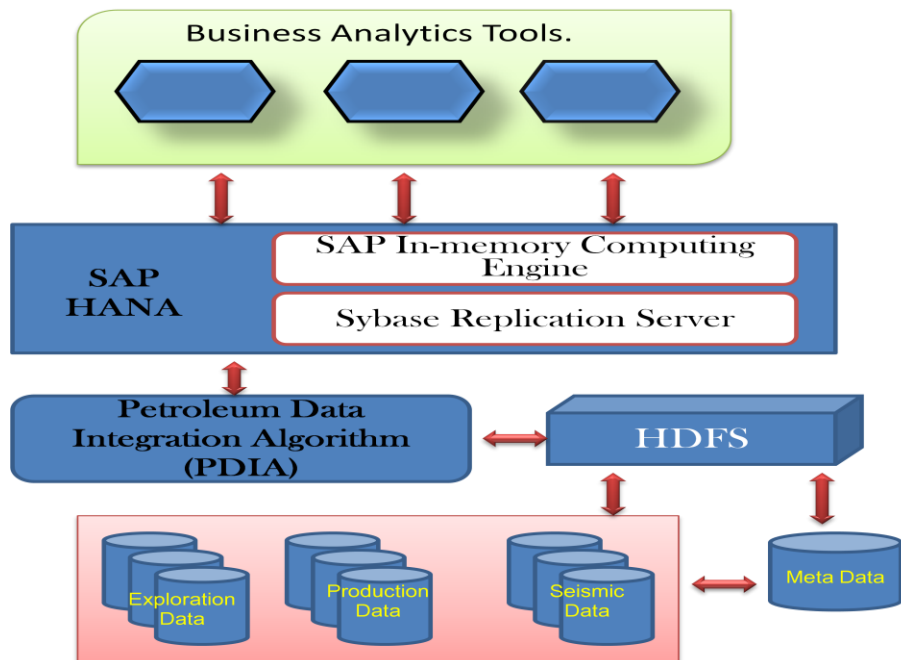


Figure 3.4. Component Diagram of HANA Engine

3.3. IDENTIFICATION & INTEGRATION ALGORITHM.

3.3.1. INTRODUCTION.

The required algorithm should not only be able to identify the proper data faster and extract it to the HANA database to make the system real-time compliant. The algorithm stated below tries to find an optimum solution for the present problem by approaching a three dimensional table concept and bring in the required datasets into the HANA database uniquely identified by their Asset-ID. The algorithm then applies different association rules to be able to find the required output in real-time scenario.

This enables a continuous real time data extraction and integration of E&P data into the in-memory HANA Database for further processing and analytics. Through this process the basic problem identified in the beginning of the research of manually locating and extracting data for different analysis would be converted to automatic, convenient and highly accurate.

3.3.2. DATA MODEL.

When we are handling petabytes of data the problem is the huge amount of time consumed in serial processing which normally may run in hours or even days. To handle the problem the algorithm should be able to implement parallel processing as the amount of data to be handled runs into petabytes normally.

To enumerate the other properties of the algorithm which are required to enhance the throughput of the system are given below.

1. Requirement of extremely high read/write operations
2. Economical scanning of data sets and subsets
3. Proficient large scale joining of datasets through one-to-one and one-to-many relationship.
4. Keep record of change of datasets over time through “time-stamping”.

3.3.3. STORAGE SYSTEM.

The storage system could be described as a light, multi-level, distributed sorted map. The map in question is indexed three ways to make an efficient and fast identification and retrieval of data. Indexing is done by row key, column key and timestamp.

I. Rows:

The row key can be defined as an atomic integer in any case of the number of column being written or read in that same row. This gives efficiency to the identification of data in spite of presence of concurrent updates to the same row.

Row keys are essentially the “Asset ID” of the exploration or production unit. It is given to be able to identify uniquely each source of generated data. In exploration and/or production whether it is a production well or exploratory well it is called an asset and given a unique ID which is called

Asset ID. This enables reduced communication with time and interaction with a limited number of systems.. The search algorithm first refers to the row key as per the request made by the client, the master indexer contains the location of the required row key mapping to get the data at a higher rate of access.

II. Column:

Column keys are assembled into sets identified as “*column genus*”, this is the basis of access control. In a column genus same type of data is stored (e.g. electric log, gama log, etc). The column genus has to be separately generated before data can be stored in it. The number of genus is kept low and less susceptible to any change during the different operations.

A column key consists of “Genus + ID”, the ID is an arbitrary integer. The genus is a name of the database type in which it is saved (e.g. Oracle, MS-SQL, MySQL, etc) and the ID is given to any particular database. The ID also enables an anchor to figure out the location of the genus.

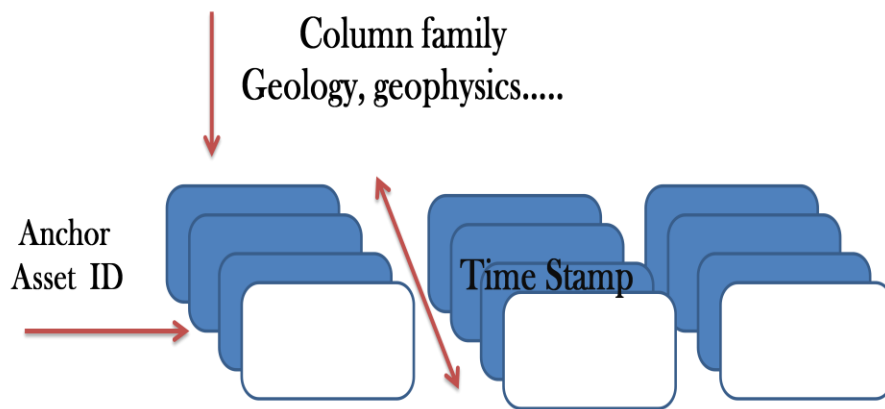


Figure 3.5. Structure Of Data Model-1

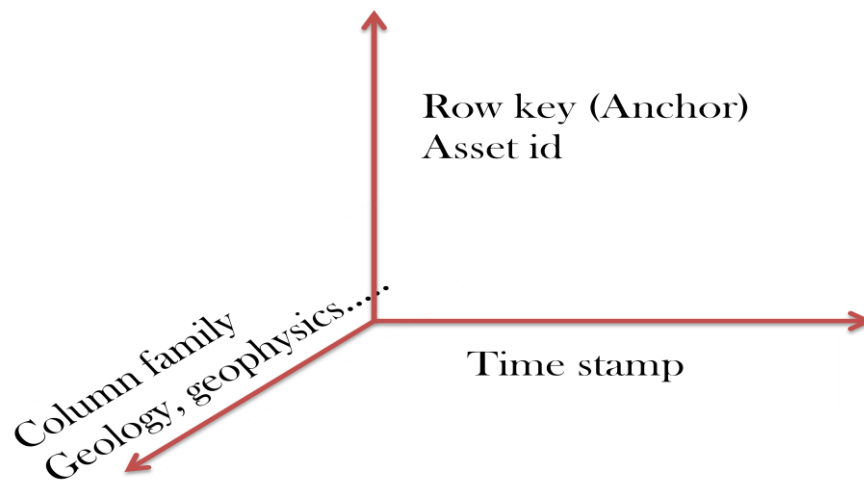


Figure 3.6. Structure Of Data Model-2

III. Timestamps:

The indexing is done by timestamps to enable multiple variants of the same data. In an electric log sample in a particular location varies with time and that is required to understand the current scenario of the reservoir and compare it with the former one so that a prediction could be made accordingly. The timestamp is generated by the synchronized system through the database. The indexing by the timestamp is done by descending order so that the most recent data is kept at the top which is read first.

3.4. DATA IDENTIFICATION AND ASSOCIATION ALGORITHM.

Algorithm 1:

// create table in HANA

1. Create Table *T = OpenTable ("Virtual table Name in HANA");
2. Row Mutation $r_1(T, \text{"asset id"})$; // Write a new anchor and delete an old anchor
3. $r_1.Set(\text{"anchor:Asset ID"}, \text{"376594"})$;
4. $r_1.Delete(\text{"anchor:12345"})$;
5. Operation op;
6. Apply (&op, & r_1);

In this algorithm the basic table is created in the HANA database and the table is populated from the source remote database.

Algorithm 2:

//Row mutation r_1 algorithm

//Merging of two partial tree for virtual organization in HANA

Input: T_l and T_r {two input partial trees }

Output: T_m {merged partial tree}

1. $k \leftarrow$ the least key in T_r
2. $h_l \leftarrow \text{GetHeight}(T_l)$
3. $h_r \leftarrow \text{GetHeight}(T_r)$
4. $h_{\min} \leftarrow \text{GetMin}(h_l, h_r)$
5. if $h_l \leq h_r$ then
6. $l_m \leftarrow$ the left most node in T_r at h_{\min}
7. add k to l_m
8. $merged \leftarrow$ merge the root of T_l and l_m
9. return $T_m \leftarrow T_r$
10. else
11. $r_m \leftarrow$ the right most node in T_l at h_{\min}
12. add k to r_m
13. $merged \leftarrow$ merge the root of T_r and r_m

14. return $T_m \leftarrow T_1$

15. endif

This algorithm performs row mutations in HANA database to create the required primary table on which further operations are performed.

Algorithm 3:

//User input required

1. Anchor Asset ID
2. D data type required
3. R Range of Time Stamp
4. C Association condition for data

These are the user inputs which are received from the user interface.

// Read Data from Table in HANA

1. Scanner scanner(T);
2. ScanStream *stream;
3. stream = scanner.FetchColumnFamily("Anchor");
4. stream SetReturnAllVersions("R");
5. scanner.Lookup("D");
6. for (; !stream Done(); stream Next()) {
7. printf("%s %s %lld %s\n",
8. scanner.Anchor(),

```
9. Stream ColumnName(),
10. Stream Timestamp (),
11. Stream Value ();

}
```

Algorithm 4:

// Searching for null value in the streamed data

1. For all streamed data
2. If $D(\text{value}) =$
3. Delete from stream r_1
4. Else
5. Send D to Procedure (C).
6. End If.

// procedure association

1. If $C \neq 0$
2. Apply C on Stream r_1
3. print Output
4. else
5. Print r_1
6. End If

3.5. WORKING OF THE ALGORITHM.

The above given algorithm basically works in two stages to form a table, in the first stage it accepts the asset-id from the user and creates a table in the HANA database by pulling in data from all the remote databases where the asset-id is common, and then the different rows are joined together to form the base table which contains all the different tables of an asset together. Next part creates a snow flake schema in the database for handling of three dimensional data as shown in figure 3.8. then further processing of data is performed according to the user input described in the algorithm.

3.6. IMPLEMENTATION OF THE PROPOSED ARCHITECTURE.

In the formal circumstances the architecture of the whole architecture is out of question for any E&P company, so the way out is to be able to modify the existing architecture in a way that neither it causes an abrupt shutdown of the day to day business nor require a huge investment for a parallel infrastructure. Thus the architecture has to be such which needs to be able to incorporate the existing hardware in such a way that it does not create any disturbance during the procedure of implementation. This could be easily possible for this architecture because of two reasons, primarily of the cloud storage which could incorporate any type of software as well as hardware which are already in action and secondly a light and agile algorithm which is able to communicate to any type of system at a very high speed with minimum failure rate.

3.6.1. THE DATA IDENTIFICATION SCENARIO.

In the industry today the different ways to associate data is by asset-id, field-id, etc which is described in figure 3.7. as given below. This leads to ambiguity in the data identification process and loss of time. The standard process could feature any one of them which could relate to any other of the components. This not only reduces ambiguity but loss of time. Under certain condition, if it is required to find the associated data with something else other than Asset-ID, it could also be achieved by the relationship with the Asset-ID. So ultimately it does not create any inadequacy in the system.

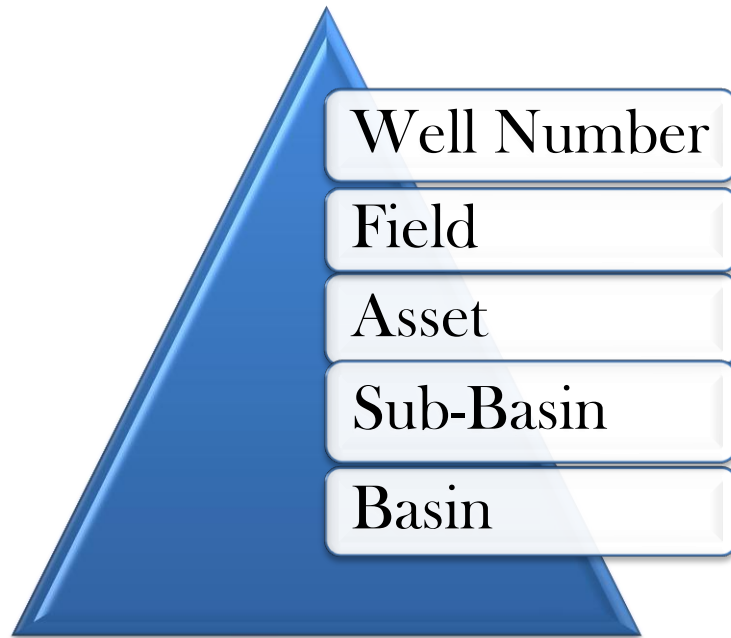


Figure 3.7. Data Relationship

3.6.2. SOFTWARE WORKING DESCRIPTION.

The developed software is a prototype to act as a technology demonstrator; it is not a fully fledged deployment package. This software demonstrates the successful working of the algorithm designed for the solution of the problem at hand. It is not in any way made to handle each and every aspects of E&P data manipulation in its current condition. Due to the constrained time and funds, developing such a package is beyond the scope of this work. But the basic structure of the software is well grounded and tested to work efficiently as expected. So any further work could easily be based on this work without any or with a minimum changes.

3.7. Testing.

In this chapter the testing of the above algorithm is described. It discusses about the setup of the testing platform, the tables, the databases as well as the forms created.

3.7.1. TESTING ARCHITECTURE.

A scaled down architecture is designed for the testing purpose. The design of which is pictographically represented in figure 3.7. In this architecture for the convenience of understanding and monitoring of data integration three separate databases are created in MS-SQL Server 2010. These three data bases contains multiple table of the same table consisting of

1. Well Table
2. Pressure & Temperature Table
3. Production Table
4. Steam Injection Table
5. Assay Table
6. Structure Table

The components of the tables are same as would have happened if the databases would have been in three different locations of drilling. The concept here is to create a scaled down version of the real scenario.

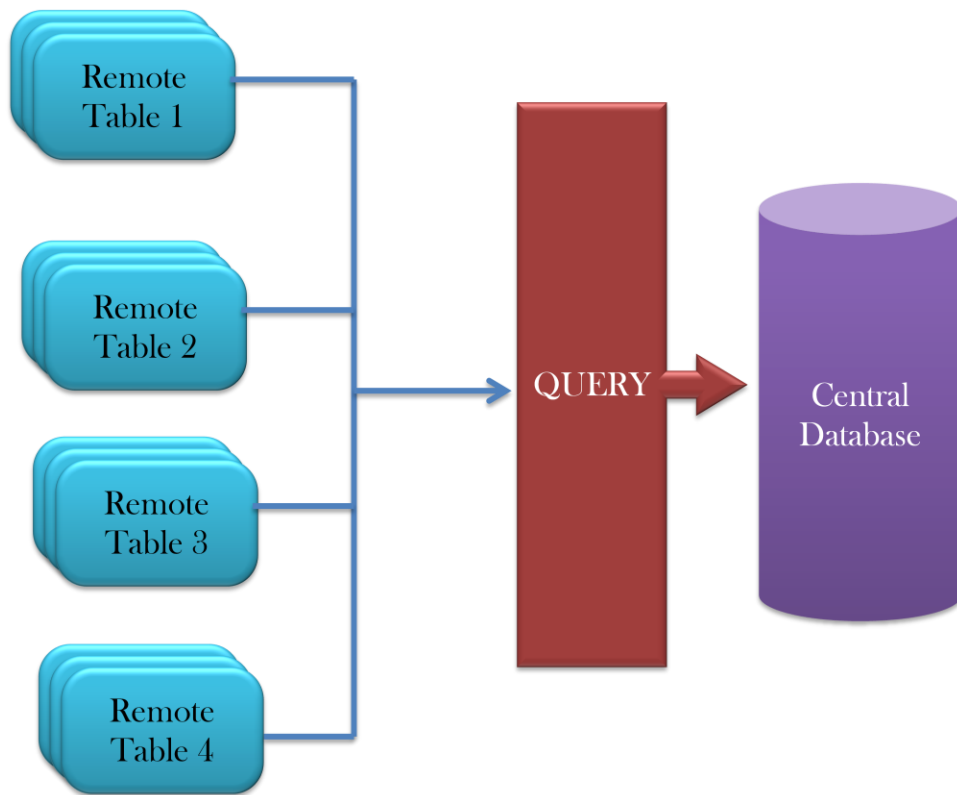


Figure 3.8. Architecture of the Setup.

3.7.2. THE SETUP.

The environment consists of two major components the databases are created in MS-SQL Server 2010 and the Human Machine Interface (HMI) is developed in Visual C++. The operating system used is Windows 7.

The fields and the data types are of the three tables are described in the following three tables. The data entered in these tables are random numbers confirming to the required data types. This enables the system to be tested correctly as far as possible without infringement of any type of national and international law. The tables are as designed below.

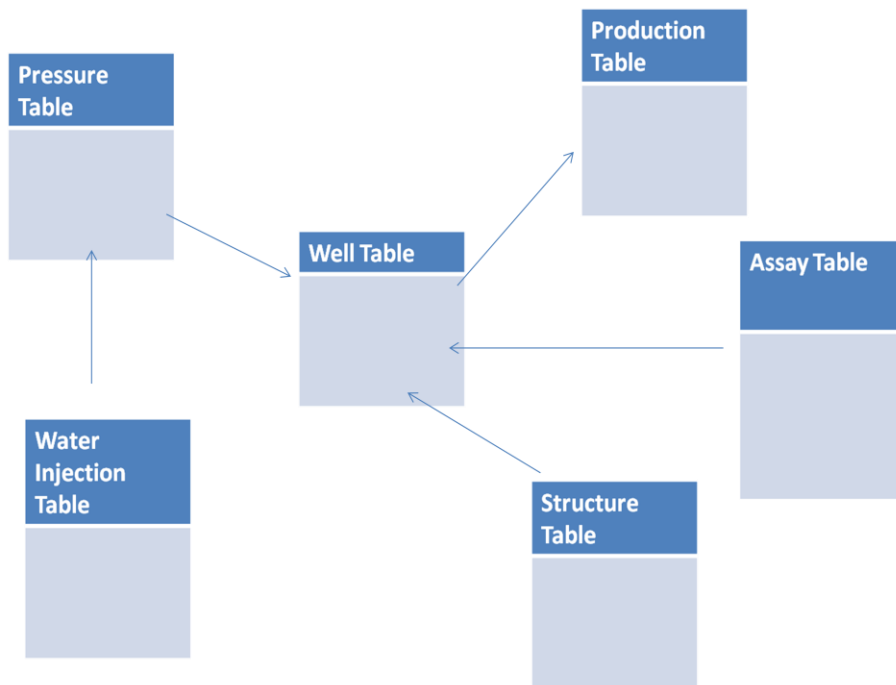


Figure 3.9. Primary Table Structure.

These are the six tables which to be created for testing purpose. The structure of these tables is described below.

3.7.3. TABLE STRUCTURE.

Fields	Data Type
Asset ID	Varchar2
Co-ordinate -X	Float
Co-ordinate -Y	Float
Co-ordinate -Z	Float
Condition	Varchar2
Type	Varchar2

Table 3.3. Asset

Fields	Data Type
Asset ID	VARCHAR2
Drill Head ID	VARCHAR2
X	FLOAT
Y	FLOAT
Z	DOUBLE
MD	INTEGER
INCL	FLOAT
AZIM	FLOAT
DX	FLOAT
DY	FLOAT
TVD	INTEGER
Calc DLS	FLOAT
Time Stamp	DATE-TIME

Table 3.4. Well

Fields	Data Type
Asset ID	VARCHAR2
Cu (ppm)	INTEGER
AU	FLOAT
Interval (m)	FLOAT
Sample ID	VARCHAR2
Time Stamp	DATE-TIME

Table 3.5. Assay

Field	Data Type
Asset ID	VARCHAR2
Type	VARCHAR2
α	INTIGER
β	INTIGER
Calc Dip	FLOAT

Calc Dim	FLOAT
Time Stamp	DATE-TIME
Null value accepted	

3.6. Structure.

Fields	Data Type
Asset ID	VARCHAR2
BWPD	Integer
Time Stamp	DATE-TIME

Table. 3.7. Steam Injection

Fields	Data types
Asset ID	VARCHAR2
well	Integer
Oil Production(BBL)	Integer
Oil Production cumulative (BBL)	Integer

Gas Production(mcf)	Integer
Gas Production cumulative (mcf)	Integer
Time Stamp	DATE-TIME

3.8. Production

Fields	Data types
Asset ID	Varchar2
Record time	Date Time
Time elapsed	Time
Pressure (Psia)	Float
temperature	Float
comments	Varchar2

Table 3.9. Temperature & Pressure

The cube structure is as follows.

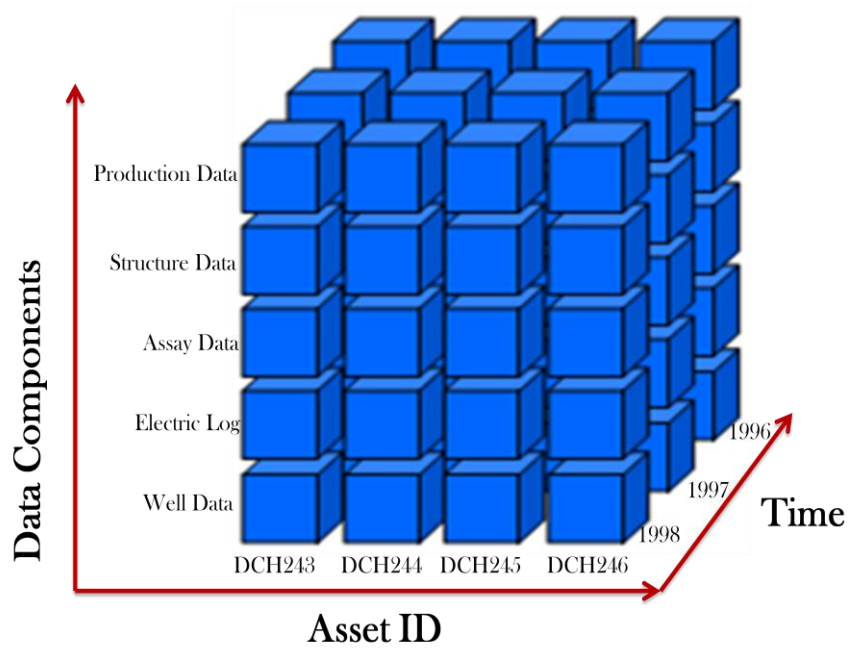


Figure 3.10. Data Cube Structure.

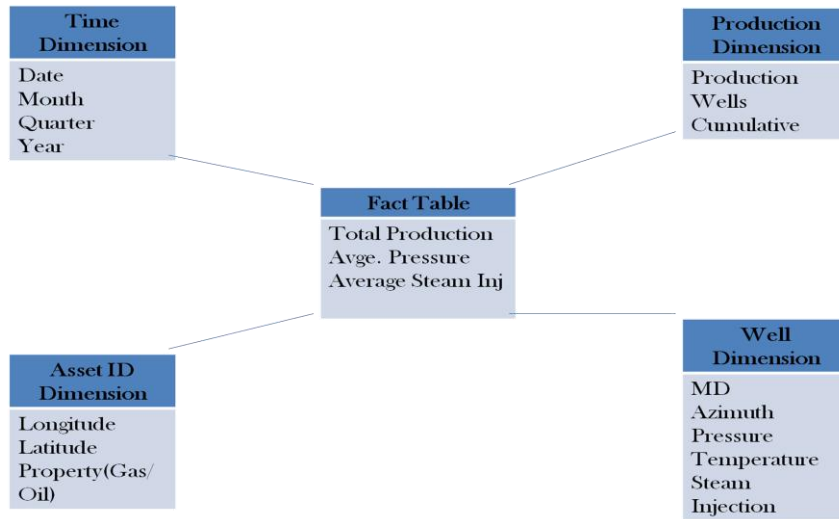


Figure 3.9. The star schema

For processing purpose this star schema is created to handle the three queries by the intended software, the queries are as follows.

The three questions to be answered by the system are

1. Maximum producing oil well.
2. Average water injection in a year.
3. Minimum fluid pressure in the last quarter.s

3.7.4. THE FRONT END.

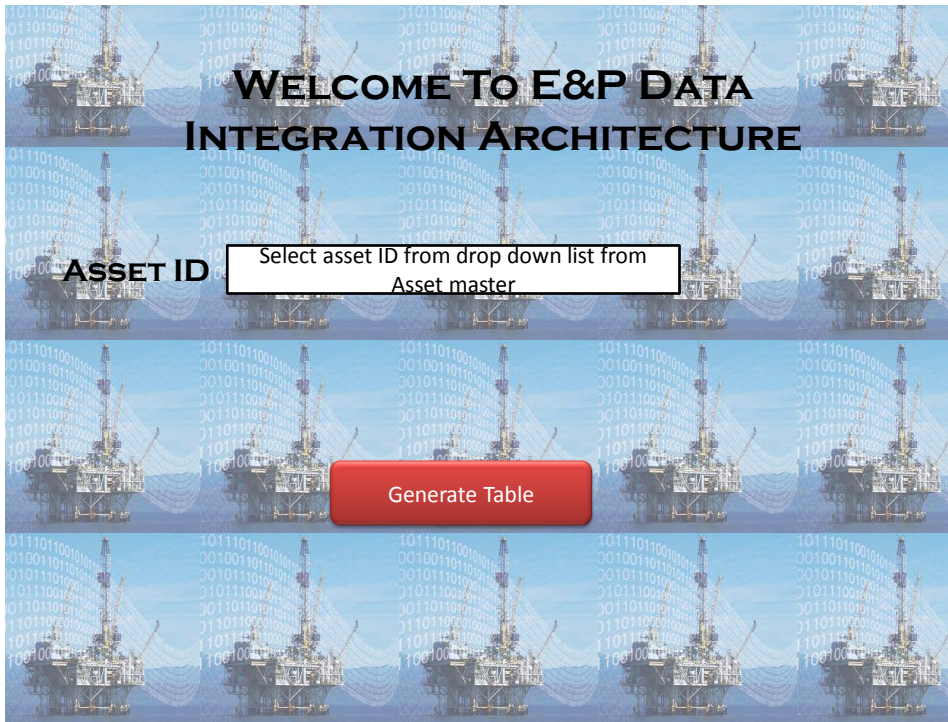


Figure 3.12. User Interface 1.

1. Query for generating list of asset_id

```
SELECT ASEST_ID FROM ASSET_MASTER
```

For the button “Generate Table”.

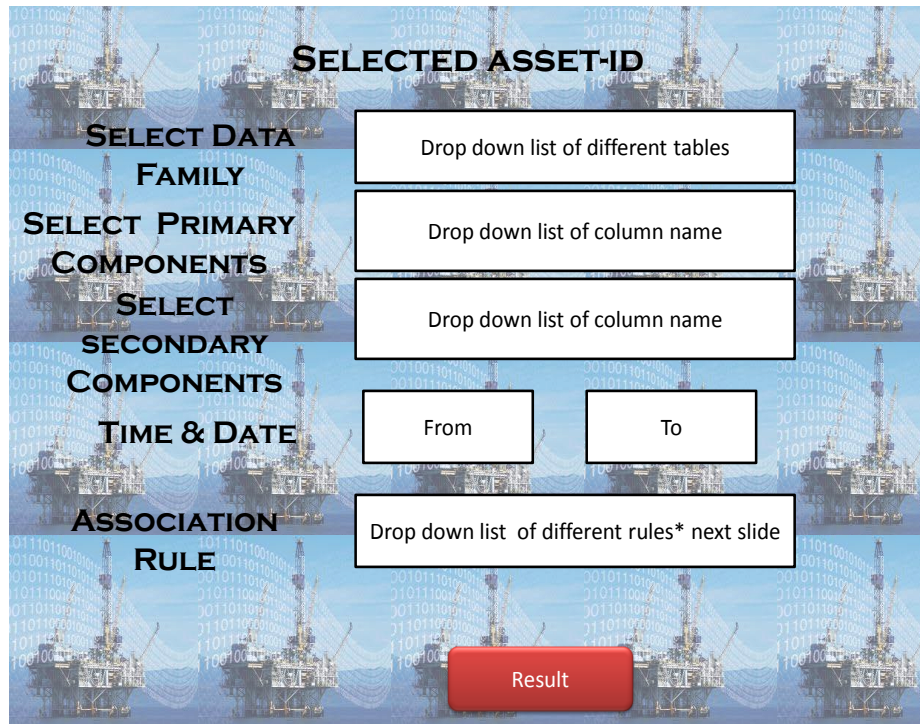


Figure 3.13. User Interface 2.

1. Query for “select data family”

----- Code Start -----

```
declare @val_to_search varchar(50), @column_name varchar(50)
Select @val_to_search = ' DDH245', @column_name='Asset_ID'
```

```
declare tbl cursor for
```

```
select table_name from information_schema.columns where
column_name=@column_name
```

```
declare @tablename varchar(200),@qstr varchar(max)
declare @datapen table(table_name varchar(200))
```

```
open tbl
fetch tbl into @tablename
while @@fetch_status=0
```

```

begin
select @qstr='select top 1 '''+@tablename+''' from
'+@tablename+' where '+ @column_name + '=''' +
@val_to_search + ''''
insert into @datapen
exec(@qstr)

```

```

fetch tbl into @tablename
end
close tbl
deallocate tbl
select * from @datapen pen

```

----- Code End -----

2. Query for Select primary & secondary component. (List of column name in a table)

```

SELECT COLUMN_NAME,* FROM
INFORMATION_SCHEMA.COLUMNS WHERE TABLE_NAME =
'selected table' ORDER BY ORDINAL_POSITION

```

3.7.4.1. ASSOCIATION RULES.

1. Minimum
2. Maximum
3. Sum
4. Average
5. Difference

3.7.4.2. QUERIES FOR ASSOCIATION

1. Query for Maximum production

```
SELECT Asset_ID, MAX(Oil_Production_Cumulative_BBL) FROM  
Production_2 group by Asset_ID
```

2. Average steam injection.

```
SELECT AVG (DISTINCT BPWD) FROM SteamInj_3
```

3. Maximum pressure

```
SELECT MAX(Pressure_PSI) FROM Pressure_Temp_2 WHERE Real_time  
>='2013-07-05 11:15:00.000' AND Real_time <='2013-07-05 12:10:00.000'.
```

4. Minimum pressure

```
SELECT MIN (Pressure_PSI) FROM Pressure_Temp_2 WHERE Real_time  
>='2013-07-05 11:15:00.000' AND Real_time <='2013-07-05 12:10:00.000'.
```

Questions to be answered

The three questions to be answered by the system are:

1. Maximum producing oil well.
2. Average water injection in a year.
3. Fall of fluid pressure in the last quarter.

3.7.5. WORKING OF THE SOFTWARE.

To begin with data is entered as per the constraints given in the databases into all the tables in the three different databases. Then the update button is executed for inclusion of all those data into the central data warehouse to the predefined specific data bases created as per the requirement.

Once the backup button is clicked it initiates another form consisting of start and stop button. This commences the Identification & Integration Algorithm; the algorithm looks into the different databases and equates the time stamp of the central data warehouse tables with the remote databases time stamp, if the time stamp of the remote database is greater than that of the central data warehouse. The algorithm generates an updated query automatically thus updating the data into the central data warehouse.

3.7.6. THE BENEFITS OF THE SYSTEM.

Here the salient features and the benefits of the designed system is enumerated below for clear understanding of the properties.

1. Streamlined architecture.
2. Utilization of existing hardware infrastructure.
3. Easy integration of existing tools.
4. Highly reduced cost of technology migration.
5. Full control of the data flow.
6. Implementation of cloud makes it highly platform independent.

7. Very high throughput by utilizing in-memory computing and optimized algorithm.

3.8. Results

Within the constrained limits of the testing the algorithm shows expected output on all the aspects of data identification and integration. The concept behind the algorithm being simple makes it highly efficient and fast. The time of query generation was set from every 4seconds to 20 seconds and each time ran flawlessly. The targeted destinations of the records in the central data warehouse were accomplished accurately with all entries. Even when the amount of data was increased the results were absolute.

CHAPTER 4

CONCLUSION

The main objective of this research work has been to identify and relegate the problems of data handling in the field of E&P sector. The varied type of data being handled in this field as well as the humongous quantity of data is an enormous problem in itself. But that is the challenge we are required to overcome. This research has tried to enumerate the different problems which are required to be detached from the data handling system, to pave the way for a real time architecture.

A lot of discussion has been done on the current market scenario of E&P sector in chapter 2 of this thesis. Other than that some even recent development worth mentioning are by Saudi Aramco[156] which is working on a conventional Data Warehouse project for handling its E&P data. The other one is by Microsoft known as Microsoft Upstream Reference Architecture Framework (MURA) [146]. , which is implementing the concept of Big-Data in the system. The basic drawback with a conventional Data Warehouse is, it has a high time lag for extraction of data from transactional systems, and the huge volume of data also augments the problem even further, so a safe conclusion could be drawn that the architecture does not support real time concepts. The Microsoft product is proprietary software which is based on the concept of big data so, it could be

safely presumed that it would be able to handle the enormous volume of data an E&P company is likely to generate. On the other hand high cost could be one of the major issues with it, not only for the software but also the hardware migration could lead to huge capital expenditure for a company.

Whereas the architecture developed in this research looks into the entire above mentioned potential problem. The architecture is far different from the conventional Data Warehousing architecture, fine tuned to make it as near as possible to Real-Time, it works with the existing infrastructure and there is no vendor locking. It also makes it a very cheap and viable option with the integration of open-source software. The SAP HANA also gives it an edge over any of the existing competitors by enabling higher throughput at an unbelievable speed of processing.

The technology present today is highly efficient as well as effective. The only requirement is to formulate a specific solution for the required problem and modify accordingly.

The technologies like ETL, Cloud Computing are perfect to create a solution for the specific problems faced by an E&P company in its IT Infrastructure. All it needs a streamlined architecture with the present efficient technologies and some innovating thinking.

Real-Time doesn't have a definite or a single definition so it entirely depends on the realization of the requirement of the business an organization is handling.

Almost every chief E&P Company follows analogous workflow process for transmitting production data en route for office. Additionally there are homogeneous array of infrastructure and applications which support are aligned with projected workflows to acquire and supply production data into enterprise systems like SAP or JD Edwards. The main objective is to augment automation and diminish time and labor costs from work processes while obtaining timely production information from the field.

The best solution for any problem is the simplest one. In the above depicted architecture the base platform of all real time data generation like SCADA system is based over a cloud platform, which in turn will provide a highly virtualized environment to enable unimpeded communication among the different heterogeneous data sources. With cloud storage platform the most inconvenient problem which has dogged illustrate E&P industry till date, of heterogeneous source databases for integration would invariably reduce the time and effort for preprocessing of data required before diverse database integration.

To be able to extract and integrate data at high speed for enabling real time data integration the algorithm needs to be simple, light as well as agile. The algorithm designed here is very simple in its approach thus it is light. It generates a search command at a regular interval of time and follows all the enlisted distributed databases each at a time, whenever it encounters a new entry in it, it copies that entry into the central data warehouse with a time and date stamp. The next time

the query visit the same database it looks for the last update date and time stamp and any entry later than that is again updated.

After implementation of the architecture with the developed algorithm the test results has shown to be highly effective in addressing the initial problem stated at the beginning. The tests have been conducted in a controlled and limited environment with emphasis on the speedy and flawless execution of the developed algorithm.

SCOPE OF FUTURE WORK

In this research the emphasis was on the identification and integration of data, but there is another aspect of quality which needs to be looked into. Any data which is to be put into a system for analytical processing needs to be “cleaned” and checked for any kind of discrepancies. Thus a lot of scope is present in this sphere for future work. The major concern in this research has been the different aspects of cloud computing as a whole, which has not been able to provide a completely reliable environment to the user as a system. A lot of concerns and apprehensions are present today with all the products available today in the market.

One of the biggest concerns in present day cloud system is the ownership of the data and its security. This is more the managerial aspect than the technicality of the problem. There is a lot of research opportunity present today in this area of cloud computing.

Another big problem the cloud computing is being ploughed by is vulnerability of the client of being over dependent on a single provider and the products they market. It limits the different options which could be profitably availed in the market.

BIBLIOGRAPHY

- 1 A. Cal` , D. Calvanese ,G. De Giacomo, and M. Lenzerini.(2002),
On the expressive power of data integration systems.
- 2 Abadi, Daniel J. (2009). Data Management in the Cloud:
Limitations and Opportunities. Bulletin of the IEEE Computer
Society Technical Committee on Data Engineering. Pages- 1-10.
- 3 Alkis Simitsis, et.al. A methodology for the conceptual modeling
of ETL processes. CEUR-WS/Vol-75.
- 4 Amrhein, D. & Willenborg, R. (2009), ‘Cloud computing for the
enterprise, Part 3: Using WebSphere Cloud Burst to create private
clouds’
- 5 Apers, P., Hevner, A., and Yao, A., “Optimization algorithms for
distributed queries” , IEEE Transactions on Software
Engineering, 9 (1), pp. 57-68, 1983
- 6 Arens, Yigal, et. al. (1993). Retrieving and Integrating Data from
Multiple Information Sources .International Journal of

Cooperative Information Systems (IJCIS),2(2):127–158.

- 7 Arens, Yigal, et.al.(1993). Retrieving and Integrating Data from Multiple Information Sources .International Journal of Cooperative Information Systems (IJCIS),2(2):127–158
- 8 Armbrust, M; Fox, A; Griffith, R; Joseph, AD; Katz, RH; Konwinski, A; Lee, G; Patterson, DA; Rabkin, A; Stoica, I & Zaharia, M (2009), ‘Above the Clouds: A Berkeley View of Cloud Computing’. Technical Report No. UCB/EECS-2009
- 9 Asanovic, K; Bodik, R; Catanzaro, BC; Gebis, JJ; Husbands, P; Keutzer, K; Patterson, DA; Plishker, WL; Shalf, J; Williams, SW; Yelick, KA (2006), ‘The Landscape of Parallel Computing Research: A View from Berkeley’. Technical Report No. UCB/EECS-2006
- 10 Aumeller, D., Do, H. H., Massmann ,S.,& Rahm, E. (2007). Schema matching with coma++, Acm sigmod (p.906-908)
- 11 Aumeller,D., Do,H.H. ,Massmann ,S., & Rahm, E. (2005). Schema matching with coma++, Acm sigmod (p.906-908).

- 12 Ayres, I.(2007) Super Crunchers: Why Thinking-by-Numbers Is the New Way to Be Smart
- 13 Babuand S. Wisdom J.,(2001)“Continuous Queries over Data Streams,” ACM SIGMOD Record, vol.30, no.3.
- 14 Bagul, S.S., Ranade, N., Sharma, A. et al (2006) A Grid based Approach for Dynamic Integration and Access of Distributed and Heterogeneous Information across an Enterprise, International Conference on Information Resources Management Association, (IRMA), 2006
- 15 Barham, P; Dragovic, B; Fraser, K; Hand, S; Harris, T; Ho, A; Neugebauer, R; Pratt, I & Warfield, A (2003), 'Xen and the Art of Virtualization', Technical report, University of Cambridge Computer Laboratory
- 16 Barr, J. (2006), ‘Amazon EC2 Beta’
- 17 Barroso, L. A.; Hoelzle, U (2009), “The Datacenter as a Computer”, Morgan and Clay pool Publishers
- 18 Batini, C., Lenzerini, M., Navathe, S.B.(1986). A comparative

analysis of, Methodologies for database schema integration,
ACM Computing, Surveys, 18(4),323-364.

19 Batini,C.,Lenzerini,M., Navathe,S.B.(1986). A comparative
analysis of, Methodologies for database schema integration,
ACM Computing, Surveys, 18(4),323-364

20 Bernstein, P., Goodman, N., Wong, E.,Reeve, C., and Rothnie, J.,
“Query processing in a system for distributed databases (SDD-1)”
, ACM Transactions on Database Systems, 6 (4), pp. 602-625,
1981

21 Bernstein,P.A., Melnik,S.,&Churchill, J.E.(2006). Incremental
schema, matching, VLDB.

22 Bernstein,P.A., Melnik,S.,&Churchill, J.E.(2006). Incremental
schema, matching, VLDB.

23 Brandic, I; Music, D; Leitner, P; Dustdar, S (2009), ‘VieSLAF
Framework: Enabling Adaptive and Versatile SLA-
Management’. Gecon09. In conjunction with Euro-Par 2009, 25-
28 August 2009, Delft, The Netherlands

- 24 Buyya, R., and Venugopal, S., July 2005 “A gentle introduction to grid computing and technologies” , Computer Society of India communications,
- 25 Carey, M.J., Haas, L. M., Schwarz, P.M., Arya, M., Cody, W.F., Fagin,R., Thomas,J., H,J., and Wimmers,E.L. (1995). Towards heterogeneous multimedia information systems: The garlic approach. In RIDE-DOM, pages124-131.
- 26 Catteddu, D; Hogben, G eds. (2009), ‘Cloud Computing - Benefits, risks and recommendations for information security’, European Network and Information Security Agency (ENISA)
- 27 Challengers (2009), ‘Final Research Agenda on Core and Forward Looking Technologies’
- 28 Chappell, D. (2008), ‘Introducing the Azure Services Platform’
- 29 Chatterjee Paresh, (2010) Tech&Trend, PCQUEST February edition .pp.33-35
- 30 Chebib Soubhi,(2001), Data: An asset of critical importance, Oil Review Middleeast.

- 31 Chong, F; Carraro, G & Wolter, R (2006), 'Multi-Tenant Data Architecture'
- 32 Claiborne Courtney, et. al. (2002)Data Integration: The Key to Effective Decisions.
- 33 Clarke, G (2005), 'Open source taking over Europe - We just don't know it'
- 34 Cohen William w.,(2000) Data Integration Using Similarity Joins and, A Word-Based Information Representation, Language, AT&T Labs—Research , Shannon Laboratory.
- 35 Cohen, W.W.(1998). Integration of heterogeneous databases without common domains, Using queries based on textual similarity. Proceedings of ACM SIGMOD-98.
- 36 Crawford. Mark L., Brulé Michael, Charalambous Yanni, and Crawley Charles., (2009), SPE-JPT November edition pg 48-53
- 37 d'Aquin, M., Doran, P., Motta, E. ,and Tamma,V. A.M.(2007) .Towards a parametric ontology modularization framework based on graph transformation. In WoMO.

- 38 DeCandia, G.; Hastorun, D.; Jampani, M.; Kakulapati, G.; Lakshman, A.; Pilchin, A.; Sivasubramanian, S.; Vosshall, P. & Vogels, W. (2007), 'Dynamo: Amazon's Highly Available Key-value Store'
- 39 DG Information Society and Media - Directorate for Converged Networks and Service (2009), 'Towards A European Software Strategy - Report Of An Industry Expert Group'
- 40 European Commission (2007), 'building the e-Infrastructure: Computer and network infrastructures for research and education in Europe. A pocket guide to the activities of the Unit GÉANT & e-Infrastructure'
- 41 European Commission, Renewable Energies Unit (2008), 'Code of Conduct on Data Centers Energy Efficiency, Version 1.0'
- 42 Exner, V., Hirsch - Homann, M., Gruissem, and Wilhelm; Hennig,L. (2008). Plant db –a versatile database form an aging plant research. Plant Meth.
- 43 F. J. et al.(2002) An efficient region-based image retrieval framework. In ACM Multimedia Proceedings, Juan-les-Pins,

France.

- 44 Fagin, R., Haas, L. M., Herneandez, M.A., Miller,R.J., Popa,L., and Velegrakis,Y. (2009).Clio: Schema mapping creation and data exchange. In *Conceptual Modeling: Foundations and Applications*, pages198-236.
- 45 Fan, X; Weber, WD & Barroso, LA (2007), ‘Power Provisioning for a Warehouse-sized Computer’. Proceedings of the 34th International Symposium on Computer Architecture in San Diego, CA. Association for Computing Machinery, ISCA '07
- 46 Färber, Franz. Cha, Sang Kyun. (2011). *SAP HANA Database-Data Management for Modern Business Applications*. SIGMOD Record, (Vol.40,No.4) ,pages 45-51.
- 47 Fegaras, L., “A new heuristic for optimizing large queries” , Proc. of DEXA 98, pp. 726-735, 1998
- 48 Fellows, W. (2009), ‘The State of Play: Grid, Utility, Cloud’
- 49 Foster, I (1998), ‘The Grid: Blueprint for a New Computing Infrastructure’ , Morgan Kaufmann Publishers

- 50 Foster, I. (2008), 'Cloud, Grid, what's in a name?'
- 51 Friedman, M., Levy, A., and Millstein, T.(1999).Navigational plans for data integration . In Proceedings of the National Conference on Artificial Intelligence (AAAI),pages 67-73.AAAI Press/The MITPress.
- 52 G. Mclachlan and K. E. Basford. Mixture Models. Marcel Dekker, Inc., Basel, NY, 1988.
- 53 Gal,A.(2006).Managing uncertainty in schema matching with top-k schema, mappings, JoDS, 90-114
- 54 Golden, B (2009), 'The Cloud as Innovation Platform: Early Examples'
- 55 Golden; B. (2009), 'Capex vs. Opex: Most People Miss the Point About Cloud Economics'
- 56 Gounaris, A., Sakellariou, R., Paton, N. W., and Fernandes, A.A.A., "A novel approach to resource scheduling for parallel query processing on computational grids" , Distributed parallel databases, 19, pp. 87-106, 2006

- 57 Graefe, G., 1990 “Encapsulation of parallelism in the volcano query processing system” , Proc. of the ACM SIGMOD Conf. on Management of Data, Atlantic City, NJ, USA, pp. 102-111.
- 58 Graefe, G., 1993 “Query evaluation techniques for large databases” , ACM Computing Surveys, 25 (2)
- 59 Grau,B.C., Horrocks,I.,Kazakov,Y., and Sattler,U. (2007). Just the right amount: Extracting modules from ontologies. In Proceedings of WWW-2007: the 16th International World Wide Web Conference, Ban ,Alberta,Canada,May8-12.
- 60 Halevy, A. Y., Rajaraman, A., and Ordille, J. J.(2006). Data integration: The teen age years. In Dayal, U., Whang,K.-Y., Lomet, D.B., Alonso, G.,Lohman, G.M., Kersten, M.L., Cha,S.K., and Kim, Y.-K., editors, VLDB, pages 9-16.ACM.
- 61 Hansen, M., Madnick, S.E., and Siegel, M.(2003). Data integration using web services .In Proceedings of the VLDB2002 Workshop EEXTT and CAiSE 2002 Workshop DTWeb on Efficiency and Effectiveness of XML Tools and Techniques and Data Integration over the Web-Revised Papers, pages 165-182,

London, UK. Springer-Verlag.

- 62 Haraty, R.A., and Fany, R.C., (2001) "Query acceleration in distributed database systems", *Revista Comlombiana de Computación*, 2 (1), pp. 19-34,
- 63 Harris, D. (2008), 'Grid vs. Cloud vs. What Really Matters'
- 64 Hernandez, m. and Stolfo,s.(1995). The merge/purge problem for large databases. *Proceedings of the 1995 ACM SIGMOD*
- 65 Hurson, Ali R.and Bright,M.W.(1991). Multi database Systems: An Advanced Concept in Handling Distributed Data. *AdvancesinComputers*,32:149–200
- 66 Ioannidis, Y., "Query optimization", *ACM Computing Surveys*, 28 (1), pp. 121-123, 1996
- 67 J. Huang and S. R. K. et al. Image indexing using color correlograms. In *IEEE Int'l Conf. Computer Vision and Pattern Recognition Proceedings*, Puerto Rico,1997.
- 68 J. Rissanen. 1989 *Stochastic Complexity in Statistical Inquiry*. World Scientific,.

- 69 J.Linand A.O.Mendelzon, 1998, Merging databases under constraints. International Journal of cooperative Information Systems,7(1):55–76
- 70 Jimenez-Ruiz, E., Grau, B.C., Sattler,U., Schneider,T., and Llavori, R.B.(2008). Safe and economic re-use of ontologies: A logic-based methodology and tool support. In ESWC, pages 185-199.
- 71 Johnson, K. (2009). Hardware, Software Innovations. The American Oil & Gas Reporter .
- 72 Kincaid, J (2009), ‘T-Mobile Sidekick Disaster: Danger’s Servers Crashed, And They Don’t Have A Backup’
- 73 Kossmann, D., and Stocker, K., (2000)“ Iterative dynamic programming: A new class of query optimization algorithms” , ACM Transactions on Database Systems, 25 (1), pp. 43-82,
- 74 Krishnamoorthy, S., (2007) Integrated Distributed Query Processor for the Data Grids, IADIS International Conference on WWW/Internet 2007, Oct 5-8, Vila Real, Portugal (accepted for publication)

- 75 Labrinidis, A., Roussopoulos, N., , 2004 “ Exploring the tradeoff between performance and data freshness in database-driven web servers” , The VLDB Journal, 13 (3), pp. 204-255, 2004
- 76 Lazcorreta Enrique, Botella Federico, (2008). Towards personalized recommendation by two step modified Apriori data mining algorithm. Expert Systems with Applications 35. Science Direct. Pages 1422–1429.
- 77 Lenzerini, M. (2002).Data integration :A theoretical perspective. In Popa, L., editor, PODS,pages233-246.ACM
- 78 Leung, AW; Pasupathy, S; Goodson, G & Miller, EL (2008), ‘Measurement and analysis of large-scale network file system workloads’, ATC'08: USENIX 2008 Annual Technical Conference on Annual Technical Conference, USENIX Association, p. 213-226
- 79 Levy, A.Y.(1998).The information manifold approach to data integration. IEEE Intelligent Systems,13:12-16.
- 80 Li, C. ,Yerneni, R. ,Vassalos, V. ,Garcia -Molina, H., Papakonstantinou, Y. ,Ullman,J. D., and Valiveti,M.

(1998).Capability based mediation in tsimmis. In SIGMOD Conference, pages564-566.

- 81 Lin, S., , 1965 “Computer solutions of the traveling salesman problem” , Bell System Technical Journal, 44, pp. 2245-2269
- 82 Liu, C., Chen, H., 1996 “A hash partitioning strategy for distributed query processing” , Proc. of 5th Intl. Conf. on Extending Database Technology: Advances in Database Technology, LNCS Vol. 1057, pp. 373-387
- 83 Liu, L., Pu, C.,, and Richine, K., , 1998 “Distributed query scheduling service: An architecture and its implementation” , Intl. Journal of Cooperative Information Systems, 7 (2 & 3)
- 84 M. H. Dunham. Data Mining, Introductory and Advanced Topics. Prentice Hall, Upper Saddle River, NJ, 2002.
- 85 Malis, A. (1993), ‘Routing over Large Clouds (ROLC) Charter’, part of the 32nd IETF meeting minutes’
- 86 Massó, J (2009), ‘Stormy Weather (Cloud & SaaS)’. INES General Assembly Keynotes

- 87 McLeod D , Heimbigner D (1980), A federated architecture for database systems. dl.acm.org
- 88 Mell, P & Grance, T (2009), ‘National Institute of Standards and Technology, Information Technology Laboratory’
- 89 Members of EGEE-II (2008), ‘An egee comparative study: Grids and clouds evolution or revolution. Technical report, Enabling Grids for E-science Project’
- 90 Michael Brulé,et.al. Reducing the “Data Commute”Heightens E&P Productivity. SPE journal. Sep-2009,pp-48-53.
- 91 Mimno, Myers & Holum. (2007)“Selection of an ETL tool”.
- 92 Nesime Tatbul,(2010) Streaming Data Integration: Challenges and Opportunities, ETH Zurich, Switzerland.
- 93 New York Times (2001), ‘Internet Critic Takes on Microsoft’
- 94 New York: Random House. Brulé, M. Technology Opens New IT Capabilities .American Oil & Gas Reporter, November. 2008.

- 95 Next Generation GRIDs Expert Group (2006), ‘Future for European Grids: GRIDs and Service Oriented Knowledge Utilities - Next Generation GRIDs Expert Group Report 3’.
- 96 Next Generation GRIDs Expert Group (2003), ‘Next Generation GRIDs: European Grid Research 2005-2010’,
- 97 Noy, N.F. and Musen, M.A.(2004). Specifying ontology views by traversal. In McIlraith, S.A., Plexousakis,D., and van Harmelen, F., editors, International Semantic Web Conference, volume 3298of Lecture Notes in Computer Science, pages 713-725. Springer
- 98 Or, I. (1976) “ Traveling Salesman-Type Combinatorial Problems and their Relation to the Logistics of Regional Blood Banking” , Ph.D. thesis, Northwestern University, Evanston, Illinois
- 99 Ouksel, ArisM.andSheth, AmitP. (1999). Semantic Interoperability in Global Information Systems:ABrief Introduction to the Research Area and the Special Section. SIGMOD, Record, 28(1):5–12.
- 100 P Shvaiko, Euzenat J – 2005, A survey of schema-based

matching approaches Journal on Data Semantics IV,

- 101 Paige Leavitt. (2002, October) Applying Knowledge Management to Oil and Gas Industry Challenges
- 102 Palmisano, I., Tamma, V., Payne,T., and Doran, P. (2009). Task oriented evaluation of module extraction techniques. In 8th International Semantic Web Conference (ISWC2009), volume 5823 of LNCS, pages130-145. Springer.
- 103 Panos Vassiliadis, et.al.(2009) A Taxonomy of ETL Activities, Proceeding of the ACM twelfth international workshop on Data warehousing and OLAP. china. Pp-25-32.
- 104 Petry, A (2007), “Design and Implementation of a Xen-Based Execution Environment”
- 105 Rahul Sah, High Performance Analytics, PCQUEST February edition 2010.pp-22-26.
- 106 Rebecca Somers “Achieving Enterprise GIS”, , January 2005, Geospatial Solutions
- 107 Richardson, L. and Ruby, S.(2007). Restful web services.

O'Reilly.

- 108 Right Scale Inc. (2009), 'Right Scale Cloud Management Features'
- 109 Rob Karel (2007) "The Forrester Wave: Enterprise ETL,
- 110 Rochwerger, R; Caceres, J; Montero, RS; Breitgand, D; Elmroth, E; Galis, A; Levy, E; Llorente, IM; Nagin, K & Wolfsthal, Y (2009), 'The RESERVOIR Model and Architecture for Open Federated Cloud Computing'. IBM Systems Journal, September 09
- 111 S. Abiteboul, D. Quass , J.McHugh, J. Widom, and J.L. Wiener.(1997) The Lore query language for semi structured data. Int.J.onDigitalLibraries,1(1):68–88.
- 112 S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and R. Harshman. Indexing by latent semantic analysis. Journal of American Sociation of Information Science, 41: 1990.pp-391–407,
- 113 S.Abiteboul, D.Quass ,J.McHugh, J.Widom,and J.L.Wiener., 1997, The Lore query language for semi structured data. Int.J. on

Digital Libraries, 1(1): 68–88

- 114 S.Babuand J.Wisdom, (2001) “Continuous Queries over Data Streams,” ACM SIGMOD Record, vol.30, no.3,
- 115 Saleem Khalid , Bellahsene Zohra,(2007) Large Scale Automatic Schema Matching, Category of submission: Survey Paper, University Montpellier.
- 116 SAP (2011). Introduction to SAP HANA for Developers
- A Pocketbook of Tutorials. Version 2.0.
- 117 Schubert, L; Kipp, A; & Wesner, S (2009), ‘Above the Clouds: From Grids to Resource Fabrics’. In G. Tselentis, J. Domingue, A. Galis, A. Gavras, D. Hausheer, S. Krco, et al., Towards the Future Internet - A European Research Perspective (pp. 238 - 249). Amsterdam: IOS Press.
- 118 Seidenberg ,J. and Rector, A.L. (2006).Web ontology segmentation: analysis, classification and use. In WWW, pages13-22.
- 119 Selinger, P. G., Astrahan, M. M., Chamberlin, D. D., Lorie, R.

A., and Price, T.G., 1979 “Access path selection in a relational database management system” , Proc. of the 1979 ACM SGMOD Intl. Conf. on the Management of Data

- 120 Sheth, Amit P and Larson, James A. (1990). Federated Database Systems for Managing Distributed, Heterogeneous, and Autonomous Databases. *ACM Computing Surveys*, 22(3):183–236
- 121 Sheth, AmitP, Gala, Sunit K., and Navathe, Shamkant B. (1993). On Automatic Reasoning For Schema Integration. *International Journal of Intelligent and Cooperative Information Systems*,2(1):23–50
- 122 Siegler, MG (2009), ‘Animoto Is Already Cash-Flow Positive, Raises Another Round To Go To 11’
- 123 Sims, K. (2009), ‘IBM Blue Cloud Initiative Advances Enterprise Cloud Computing’
- 124 Sotomayor, B; Montero, RS; Llorente, IM & Foster, I (2009), ‘An Open Source Solution for Virtual Infrastructure Management in Private and Hybrid Clouds’. *IEEE Internet Computing*,

Special Issue on Cloud Computing, October 09

- 125 Sutter, H. (2005), 'The Free Lunch Is Over: A Fundamental Turn Toward Concurrency in Software', in Dr. Dobb's Journal, 30(3).
- 126 The NESSI-Grid Project (2008), 'Grid Vision and Strategic Research Agenda'
- 127 The ROGTEC Team. (2009, September 1). Oil & Gas News. Retrieved April 27, 2010,
- 128 Toffler, A. (1980), 'The Third Wave', Pan Books
- 129 Tracy Jenee (2005) "Herding Cats! GIS Coordination Efforts in an Enterprise System", 'Moy, Arkansas Game & Fish Commission, ArcUser Online October - December
- 130 Truffle Capital (2007), 'Truffle Capital: European Commission Recognizes the Need for a "European Strategy for Software" - Commenting on the 2007 Truffle 100 Europe, Viviane Reding Calls on Europe to Develop a Leadership Position in
- 131 Truong, H.-L. And Dustdar, S.(2009).On analyzing and specifying concerns for data as a service. In Proc of IEEE Asian-

Pacific Service Computing Conference.

- 132 Vambenepe, W (2009), “Reality check on Cloud portability” - available at <http://stage.vambenepe.com/archives/684>
- 133 Vaquero, L. M.; Rodero-Merino, L.; Caceres, J. & Lindner, M. (2009), ‘A break in the clouds: towards a cloud definition’, SIGCOMM Comput. Commun. Rev. 39(1), 50—55
- 134 W Guédria, Z Bellahsene, M Roche (2007)- A flexible approach based on the user preferences for schema matching Ieee rcis, - lirmm.fr
- 135 W Su, J Wang, F Lochovsky – (2006) Holistic schema matching for web query interfaces in Database Technology-EDBT , Springer
- 136 Wang, Z., Wang, K., Topor, R.W., and Pan, J.Z.(2008). Forgetting concepts in dl-lite. In ESWC, pages 245-257.
- 137 Wayne Eckerson and Colin White.“Evaluating ETL and Data integration Platforms”
- 138 Wayner, P (2008), ‘Cloud versus cloud: A guided tour of

Amazon, Google, AppNexus, and GoGrid' - available at [http://www .infoworld.com/d/cloud-computing/cloud-versus-cloud-guided-tour-amazon- google-appnexus-and-gogrid-122?page=0,0](http://www.infoworld.com/d/cloud-computing/cloud-versus-cloud-guided-tour-amazon-google-appnexus-and-gogrid-122?page=0,0)

139 Wiederhold, G.(1992). Mediators in the architecture of future information systems. *IEEEComputer*,25 (3):38-49.

140 Wikinomics, 'The Prosumers'- available at

141 Yu, C. T., Chang, C., and Chang, Y., 1982 " Two surprising results in processing simple queries in distributed databases" , *Proc. of 6th IEEE Intl. Computer Software and Applications Conference*, pp. 377-384

142 Yves de Montcheuil. "ETL Engine or Code Generation"

143 Zhu, F., Turner, M., Kotsiopoulos, I., Bennett, K., Russell, M., Budgen, D., Brereton, P., Keane, J., Layzell, P., Rigby, M. ,and Xu,J. (2004). Dynamic data integration using web services. In *ICWS'04: Proceedings of the IEEE International Conference on Web Services*, page 262, Washington, DC, USA. IEEE Computer Society.

- 144 Ziegler Patrick, Dittrich Klaus R.(2007), Data Integration - Problems, Approaches, and Perspectives, Database Technology Research Group, Department of Informatics, University of Zurich.
- 145 Zimory GmbH (2009), 'Zimory Enterprise Cloud – Whitepaper'
- 146 <http://aka.ms/ictt9o>
- 147 <http://aws.amazon.com/s3/>.
- 148 http://en.wikipedia.org/wiki/Cloud_computing
- 149 [http://en.wikipedia.org/wiki/John_McCarthy_\(computer_scientist\)](http://en.wikipedia.org/wiki/John_McCarthy_(computer_scientist))
- 150 http://en.wikipedia.org/wiki/Geophysical_survey
- 151 http://en.wikipedia.org/wiki/Reflection_seismology
- 152 http://en.wikipedia.org/wiki/Well_logging
- 153 <http://www.analytics-magazine.org/november-december->

2011/695-how-big-data-is-changing-the-oil-a-gas-industry

154 <http://www.ft.com/cms/s/0/2400cf6a-07e3-11e2-9df2-00144feabdc0.html> # axzz2PrBtOMwM

155 http://www.gisdevelopment.net/application/business/mi08_172.htm.

156 http://www.iaeng.org/publication/WCE2007/WCE2007_pp553-

157 <http://www.ogsadai.org.uk>

158 <http://www.opencloudmanifesto.org/Open%20Cloud%20Manifesto.pdf>

159 [http://www.safe.com/technology/spatial ETL/ overview.php](http://www.safe.com/technology/spatial%20ETL/overview.php)

160 [http://www.socialtext.net/wikinomics /index.cgi? the_prosumers](http://www.socialtext.net/wikinomics/index.cgi?the_prosumers)

161 <http://www.streambase.com/>

BRIEF BIODATA OF THE AUTHOR

Name : **Kingshuk Srivastava**
Sex : Male
Marital Status : Married
Date of Birth : 24th September 1978
Father's Name : Dr. K. P. Srivastava
Address : 11/20, B-Road, Bamungachi, Salkia, Howrah. West
Bengal. Pin:-711106. Mob:-919568533349
Present Affiliation : University Of Petroleum & Energy Studies.

Academic Qualification :

Education	Board/ Univ.	College/school	Year of completion	Score	Remarks
M.Tech	UPES	UPES,	2009	CGPA	1 st Rank in

(Petro Informatics)		Gurgaon		3.54 / 4.00	University
M.Sc.- Physics (Electronics)	L.N.M.U.	C.M.Science College Darbhanga	2006	69.75%	6 th Rank in University