



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

# IJIRSET

International Journal of Innovative Research in  
**SCIENCE | ENGINEERING | TECHNOLOGY**



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 10, Issue 10, October 2021

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

Impact Factor: 7.569

# A Case Study on Traditional Vs Cloud EDW

Yedidi Naga Shiva Harish

Gayathri Vidya Parishad, College of Engineering, Visakhapatnam, India

**ABSTRACT:** Data warehousing over the cloud environment has been studied and from the study it was understood that the data warehouse system over the cloud environment has advantages such as flexibility, scalability, dependability and reduced costs, without compromising on the security and compliance standards of the on-premise data warehousing.

## I. INTRODUCTION

Currently major business Organization's Enterprise Data Warehouses are on Teradata Intelliflex System (hardware platform for demanding data analytics) which is on-Premise service model. Keeping in view of the operational cost expenditure, there is a necessity to explore alternative avenues to reduce operational cost without compromising the database performance and user experience.

This document elaborates different cloud Data warehouse systems considered for alternative Analytics platforms.

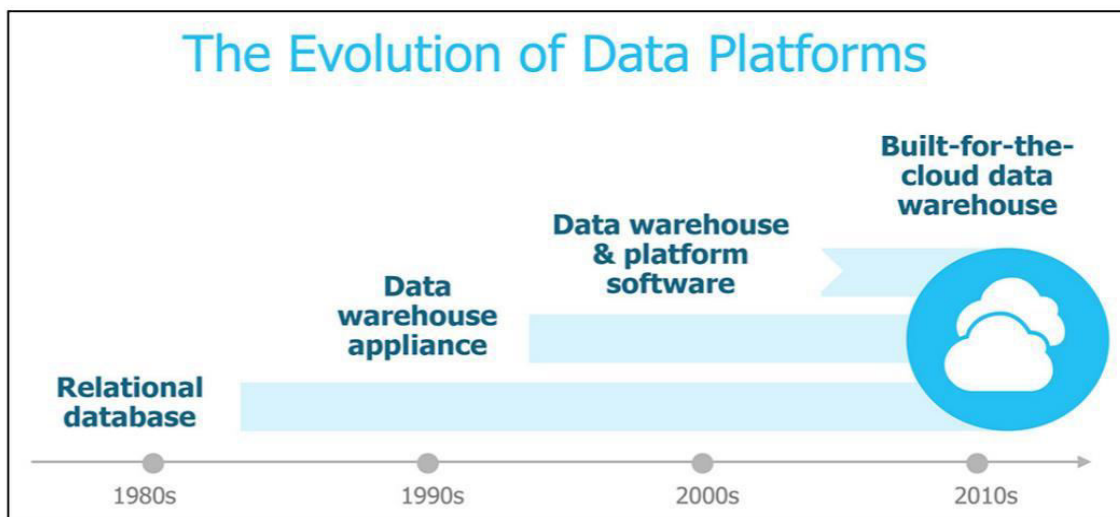
### 2. Background and Basic Terminology

#### a). What Is a Data Warehouse?

A data warehouse is a computer system dedicated to storing and analyzing data to reveal trends, patterns, and correlations that provide information and insight. Traditionally, organizations have used data warehouses to store and integrate data collected from their internal sources (usually transactional databases), including marketing, sales, production, and finance.

#### b). The Evolution of Data Warehousing

Data warehousing has evolved over four decades



#### c). Issues / limitations of the traditional data warehouse

Major setbacks of the traditional database are

- Availability
- Capacity estimation and forecasting to address exponential growth
- Performance (I/O) and user queries
- System availability and spikes
- system resource maintenance and Operational expenses (upgrades/maintenance and patching)
- Licensing costs



- Backup and DBA activity.

Cloud operating platform is emerging as alternative methodology to overcome the limitations of the traditional data warehousing.

### 3. Introduction to cloud data warehousing

Cloud data warehousing is a cost-effective way for companies to take advantage of the latest technology and architecture without the huge upfront cost to purchase, install, and configure the required hardware, software, and infrastructure. The various cloud data warehousing options are generally grouped into the following three categories

#### a). Software as a Service (SaaS)

This model allows the customers to use the entire application of the providers running on the cloud environment.

#### b). Platform as a Service (PaaS)

This model deploys a layer of software or development environment on to the cloud infrastructure and offered as a service. Higher levels of service can be built over this layer. To meet the manageability and scalability requirements of the applications, a predetermined combination of OS and application servers are suggested by the PaaS providers.

#### c). Infrastructure as a Service (IaaS)

This model provides intrinsic storage and computing proficiencies to the customers as normalized services. The basic computing resources, storage, networking equipment, etc. are shared to customers to handle workloads.

#### d). Characteristics of Cloud Computing

Cloud is expected to have the following five characteristics that are described below.

- **On-demand self-service:** Without interacting with the service providers, a customer can independently provision computing proficiencies, such as system storage and server time.
- **Broad network access:** Competences that are available over the network are accessed through typical mechanisms that confirm its use by diverse weak or strong client platforms.
- **Resource pooling:** According to the needs of customer, they are dynamically allocated and deallocated with different physical and virtual resources. The cloud providers' computing resources are then shared to various customers' using a multi-tenant model. Due to location independence, the customer has no control or knowledge over the location of available resources.
- **Rapid elasticity:** In some cases, in order to quickly scale out and scale in, automatically the proficiencies are quickly and flexibly provisioned to the customers. The proficiencies can be purchased in any number at any time by the consumer since they are unlimitedly available.
- **Measured Service:** Using a metering facility, cloud systems can automatically control and optimize resource usage based on the type of services at some level of abstraction. The utilization of resources can be controlled, observed and reported which allow both the provider and consumer to make use of their exploited services efficiently.

#### e). Advantages of cloud computing that organizations are experiencing today

- Do more with less
- Flexible costs
- Always-on availability
- Improved mobility
- Quick Software upgrades
- Improved collaboration
- Cloud computing is more cost effective
- Expenses can be quickly reduced
- Flexible capacity
- Facilitate M&A activity
- Less environmental impact

### 4. On-premise vs Cloud data warehousing: Pros and Cons

Companies are increasingly moving towards cloud-based data warehouses instead of traditional on- premise systems. Cloud-based data warehouses differ from traditional warehouses in the following ways:

- There is no need to purchase physical hardware.
- It's quicker and cheaper to set up and scale cloud data warehouses.

In recent years, many Industries' data warehouses are moving to the cloud. The new cloud-based data warehouses do not adhere to the traditional architecture; each data warehouse offering has a unique architecture.

Table below outlines differences between on-premises and cloud platforms

Sl.No.	Feature	On – Premises	Cloud
1	<b>Capacity</b>	<ul style="list-style-type: none"> <li>Requires continuous monitoring and analysis on capacity and TTL (time to)</li> <li>Capacity is provisioned based on</li> <li>guessing and periodical growth.</li> <li>For wrong estimates, we ended up paying for expensive resources and</li> <li>insufficient capacity</li> <li>Data Center resource costs</li> </ul>	<ul style="list-style-type: none"> <li>Capacity Planning is Agile</li> <li>Capacity can be facilitated in few hours</li> <li>We can add/delete capacity as temporary/ permanent based on the usage needs</li> <li>No Data center costs</li> </ul>
2	<b>Agility</b>	<ul style="list-style-type: none"> <li>Infrastructure change Planning and implementation takes time</li> <li>This delays the implementation speed.</li> <li>Affects experimentation and slow down</li> </ul>	<ul style="list-style-type: none"> <li>Experimentation quickly with low cost and risk</li> <li>Spin up resources in few clicks for innovation and testing</li> <li>Meets the organization targets in very fast due to dramatic increase in agility in implementation</li> </ul>
3	<b>Scalability &amp; Elasticity</b>	<ul style="list-style-type: none"> <li>Needs more planning and budgeting to address future resource needs</li> <li>We end up paying infrastructure or data center costs for the resources that are not used</li> </ul>	<ul style="list-style-type: none"> <li>Auto scaling and elastic load balance enables</li> <li>- Scale the resources up/down quickly and easily with quick turnaround</li> <li>- No additional costs for the resources that are not used (pay as you use model)</li> </ul>
4	<b>Security and Compliance</b>	<ul style="list-style-type: none"> <li>Organization standards and security protocols may avoid security breach.</li> <li>Manage more security access layers on Customers PII and PCI data.</li> <li>Periodic security and compliance checks will helps manage tightest securities on organization sensitive data</li> <li>Security ownership lyes with the Organization</li> </ul>	<ul style="list-style-type: none"> <li>Although cloud service providers implement the best security standards and industry certifications, storing data and important files on external service providers always opens risks</li> <li>Security Ownership shared between Cloud service provider and Organization. Improper handshaking leads to risk of security breach</li> </ul>
5	<b>Availability</b>	<ul style="list-style-type: none"> <li>Need more resources and effort to monitor systems.</li> <li>During upgrades and maintenance system will not be available</li> </ul>	<ul style="list-style-type: none"> <li>Systems always accessible with fault tolerant infrastructure</li> <li>System upgrades are seamless to end users and system is always available</li> </ul>

##### 5. Comparison study on Traditional Vs Cloud Data Warehouse Systems

We proactively conducted a proof of concept and analyzed Architectural and Technical comparison between current Traditional Data warehouse systems and different Cloud data warehouses platforms. This document details out the architectural and technical comparison of these platforms.

Following cloud platforms are considered for comparison and assessment for futuristic adoption to facilitate data analytics requirements.



- Snowflake
- BigQuery (Google Cloud Platform)

Following are the considerations against which above systems and current Teradata System are evaluated.

- Architectural -> Teradata Vs Cloud Data warehouse systems
- Features -> Teradata Vs Snowflake Vs Big Query
- Technical Differences -> Teradata Vs Snowflake Vs Big Query

Please make a note that we used Teradata prototype systems is used for this comparison.

1. Architectural differences -> Teradata Vs Snowflake Vs BigQuery

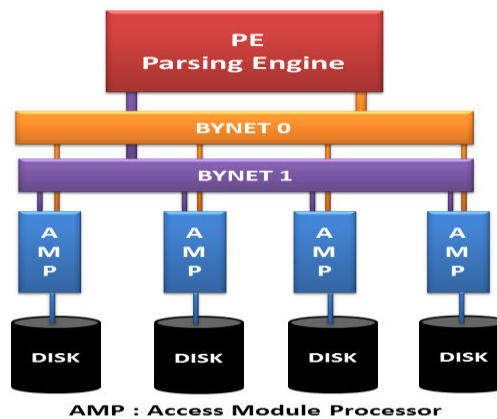
In this section, we compare the architecture of Teradata to snowflake and BigQuery.

- Teradata

The Teradata Database is a relational database management system (RDBMS) that drives an Organization data warehouse. The Teradata database is not just a software solution, the Teradata RDBMS consists of hardware AND software. The Teradata database is the first RDBMS specifically designed for Massively Parallel Processing (MPP). The Teradata Database is an open system, compliant with industry ANSI standards

Teradata Architecture and Main features

- ANSI SQL 2011 standard
- Single data store
- Linear Scalability
- Unconditional parallelism (parallel architecture)
- Ability to model the business
- Mature, parallel-aware Optimizer



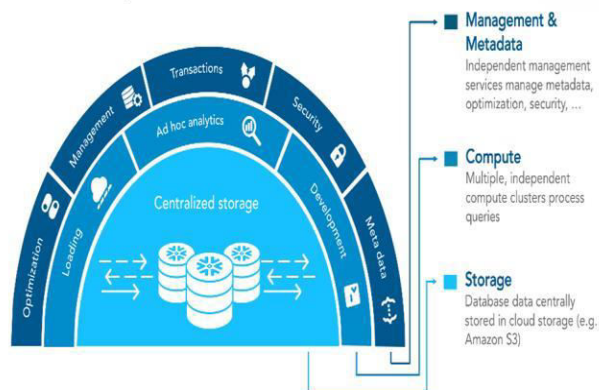
Picture Referred from: <http://www.teradatawiki.net/2013/09/Teradata-Architecture.html>

- Snowflake

Snowflake is a fully relational columnar database as an opposite to row storage databases like Oracle, MS SQL, MySQL etc. Snowflake is an analytic data warehouse provided as a Software-as-a-Service (SaaS). Snowflake provides a data warehouse that is – faster, easier to use, and far more flexible than traditional data warehouses. Snowflake has many similarities to other enterprise data warehouses, but also has additional functionality and unique capabilities.

• Snowflake Architecture and Main features

- Shared and centralized cloud storage
- Independent and flexible compute clusters
- Fully relational SQL data warehouse
- Turn-key service
- Pay what you use
- Standard interfaces for ETL and BI



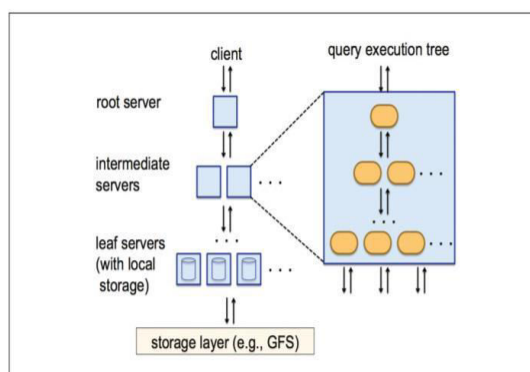
Picture Referred from [www.Snowflake.com/architecture](http://www.Snowflake.com/architecture)

### • BigQuery

BigQuery is Google's Enterprise data warehouse enabling super-fast SQL queries using the processing power of Google's Dremel architecture, providing a full-managed and cloud based interactive query service for massive data sets. This is a fast, highly scalable, cost effective and fully managed cloud data warehouse for analytics with built-in machine learning capabilities.

### BigQuery Architecture and Main features

- ANSI SQL 2011 standard
- Columnar storage
- Full or partition column scan using tens of thousands of machines over fast Google network
- Tree Architecture is used for dispatching queries and aggregating results across thousands of machines in a few seconds.
- Real-time streaming using API



Tree architecture of Dremel

Picture Referred from [cloud.google.com/igQueryTechnicalWP.pdf](http://cloud.google.com/igQueryTechnicalWP.pdf)

### 2. Feature Differences Teradata Vs Snowflake Vs Big Query

Sl.No.	Feature	Teradata	Snowflake	BigQuery
1	Initial Release	1984	2014	2010
2	Database Mode	Relational RDMS	Relational RDMS	Relational RDMS
3	Available Platforms	On Premises and Cloud	Only Cloud	Only Cloud
4	Concurrency	Managed Workloads	High concurrency	High concurrency
5	Maintenance	Needs downtime for planned and scheduled maintenances	Seamless maintenance	Seamless maintenance

6	Scalability	Linearly scalable	automatic	automatic
7	Elasticity	No	Yes	Yes
8	Parallelism	Yes	Yes	Yes
9	Data Format Support	LATIN & UNICODE	LATIN & UNICODE & Support for diverse data like JSON & PARQUET, XML, Avro	LATIN & UNICODE & Support for diverse data like JSON, Avro Parquet, ORC
10	Data operations	Supports INSERT/UPDATES /DELETES/EXPORTS	Supports INSERT/UPDATES /DELETES/EXPORTS	Append-only, batch operation SQLs are expensive
11	Data Loading	Batch Loading	Batch Loading & Streaming	Batch Loading & Streaming
12	Processing Unstructured data	No	Partially (regular expression matching on text)	Partially (regular expression matching on text)
13	Handling large results / Join large table	High Performance	Low Performance	Low Performance
14	Workload management	Yes	No	No
15	Operating Systems	LINUX/ Windows	Hosted on AWS	Hosted on GCP
16	Availability	High availability, DR built in (No backup) is not required for Intelliflex Systems. All other Teradata systems model required DR (Disaster Recovery)	High availability, DR built in (No backup)	High availability, DR built in (No backup)
17	Tools/ Utilities	Teradata has own data loading Tools BTEQ, FEXP, FLOAD, MLOAD and TPUMP	No Inhouse tools, need use vendor products or Open Source tools	Needs use other vendor products or Open Source tools
18	BI /Reporting	Supports all Analytics and Reporting tools available in the market	Tableau	Google Data Studio
19	Pricing Model	Regular	Pay-as-you-go	Pay-as-you-go
20	Machine Learning and Artificial Intelligence Analytics	No	Yes	Yes
21	Monitoring Tool	Viewpoint	No Monitoring Tool, monitor resource usage using web interface and SQL	Stackdriver

### 3. Technical Differences Teradata Vs Snowflake Vs Big Query

Sl.No.	Feature	Teradata	Snowflake	BigQuery
1	SQL	ANSI SQL	ANSI SQL	ANSI SQL
2	Stored Procedure	Yes	No	No
3	Triggers	Yes	No	No
4	Data Types	Supports Industry standard datatypes including CLOB and BLOB	Standard datatypes & semi-structured data types like VARIANT OBJECT, ARRAY. It doesn't support CLOB and BLOB	Standard datatypes & semi-structured data types like ARRAY, STRCUT Geography. It doesn't support CLOB and BLOB
5	Programming Language	Yes	Yes	Yes

	support			
6	APIs & other access methods	.NET Client API HTTP REST JDBC JMS Adapter, ODBC, OLE DB	CLI Client, JDBC ODBC	RESTful HTTP/JSON API
7	Foreign keys	Yes	No	No
8	Join Indexes	Yes	No	No
9	Secondary indexes	Yes	No	No
10	Geospatial	Yes	No	Yes
11	Compression	Yes	Yes	Yes
12	Columnar	Need define explicitly	By default	By default
13	Control Data Access	Ability to Create Roles and Views and provisioned ad secure data access	Ability to Create Roles and Views and provisioned ad secure data access	BigQuery uses Identity and Access Management (IAM) to manage access to resources
14	Partitions	Supports Single and Multilevel Partitions	micro-partitions and data clustering	Limited partitions
15	Table Types & Views	<ul style="list-style-type: none"> <li>Physical Tables</li> <li>Derived Tables</li> <li>Volatile Tables</li> <li>Global Temporary Tables</li> <li>Views</li> </ul>	<ul style="list-style-type: none"> <li>Temporary</li> <li>Transient</li> <li>Permanent</li> <li>Views</li> </ul>	<ul style="list-style-type: none"> <li>Native tables</li> <li>External tables</li> <li>Views</li> </ul>
16	User-Defined Functions	Yes	Yes	Yes
17	Materialized Views	No	Yes	No

## 6. Conclusion and recommendations

Data warehousing over the cloud environment has been studied and from the study it was understood that the data warehouse system over the cloud environment has advantages such as flexibility, scalability, dependability and reduced costs, without compromising on the security and compliance standards of the on-premise data warehousing.

## REFERENCES

1. <https://cloud.google.com/bigquery/docs/monitoring#view-dashboards> <https://docs.snowflake.net/manuals/sql-reference/data-types-unsupported.html>





**Impact Factor:  
7.569**



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 **9940 572 462**  **6381 907 438**  **ijirset@gmail.com**



[www.ijirset.com](http://www.ijirset.com)

Scan to save the contact details