

A Comparison of Cloud Data Warehouse Platforms



Sonra Intelligence Limited

GW 107, Greenway Hub

DIT Grangegorman

Dublin 7

hello@sonra.io

www.sonra.io

Contents

1. Summary.....	2
2. Cloud Platform Spending.....	3
3. The Benefits of Cloud based Analytics.....	5
3.1 Evaluation Criteria.....	6
4. Analysis of Options.....	15
4.1 Amazon Redshift.....	16
4.2 Hadoop based Data Warehouse.....	19
4.3 Microsoft Azure SQL Data Warehouse	21
4.4 Oracle Database Cloud Exadata Service.....	23
4.5 Snowflake Elastic Data Warehouse.....	25
5. About Sonra Intelligence	28
6. About the author: John Ryan	28
7. About the co-author: Uli Bethke.....	28
8. Training Big Data for Data Warehouse Professionals	29



About Sonra

Sonra Intelligence are experts in data warehouse design, implementation, and cloud migration. We provide services across the globe including [training](#) and [data architecture advisory services](#).

With offices in London and Dublin, we provide advice and guidance to blue chip clients on a range of technologies including Amazon Redshift, Snowflake, Oracle, and Hadoop.

Sonra is also developer of Flexter. [Flexter](#) is an ETL tool for XML and JSON. It automatically converts XML/JSON to text, a relational database, or Hadoop. Visit the [FAQ](#) section for common queries.

[Reach out to us](#), and see how we can help you modernise your data analytics platform.

Web: <https://www.sonra.io>

eMail: hello@sonra.io

Telephone: +353 1 5345 015

1. Summary

An ideal warehouse would be deployable within minutes on hardware which is isolated from other users, but shares access to an almost infinitely scalable data store. It would be possible to grow (or shrink) the compute resources independent of storage, with additional processing quickly deployed to cope with unexpected demands.

The solution would support an agile working method, including the ability to allocate or suspend clusters at will, and costs would be controlled in line with usage. Finally, the entire system would be simple, requiring little or no database administration or tuning.

In terms of cloud platform, Amazon Web Services (AWS) is the clear leader with a 48% market share, and two solutions, *Snowflake Elastic Data Warehouse* and *Amazon Redshift* come closest to the ideal, with both running on AWS.

While all database options perform well, Snowflake stands out as a database architected for the cloud, and has compelling benefits around simplicity and reduced skill requirements. Amazon Redshift also delivers a compelling result, although needing slightly more technical design and support.

Both Oracle and Microsoft have migrated an existing on-premises based database to the cloud, and while these provide an easier migration path for existing customers, they lose out in terms of system management and complexity.

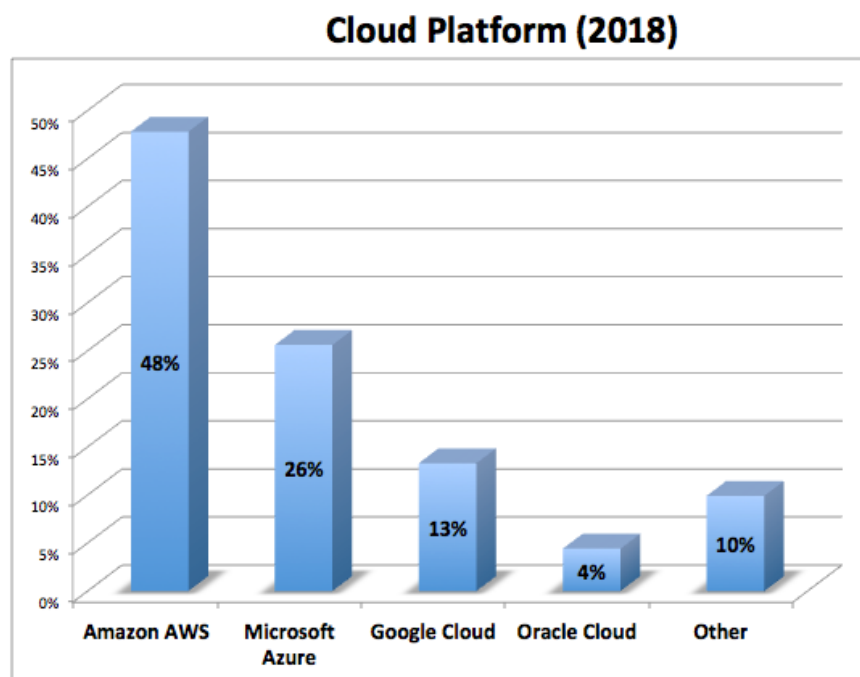
Meanwhile, despite the increasing interest in Hadoop, it has proven to be a poor option for business intelligence, and should be avoided.

The remainder of this document describes the selection criteria in detail, and compares the options available.

2. Cloud Platform Spending

In 2018 [IDC](#) estimate that half of all IT spending will be Cloud-based, rising to 70% of all software services by 2020. Meanwhile, despite the hype around Hadoop, the relational data warehouse [continues](#) to be a [critical tool](#) for any business as data volumes grow.

Another survey published by [TDWI](#) demonstrates Amazon Web Services is the market leader in cloud based services with nearly half of respondents deployed on AWS.

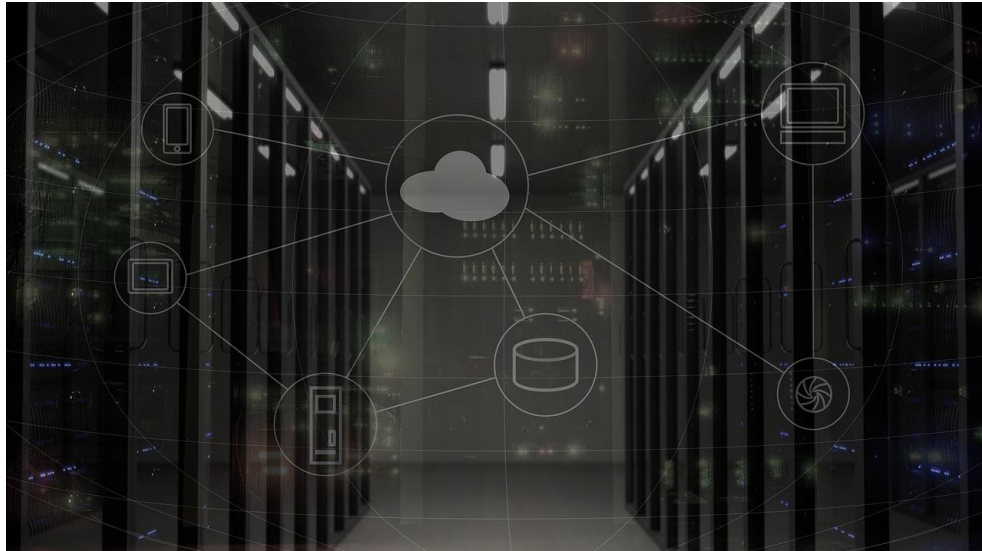


The survey found the main drivers behind cloud deployments were cited as cost reduction and control (46%) followed by agility (31%) and the ability to scale (27%).

In this paper we'll discuss the main benefits of Data Warehousing in the cloud, and outline and then explain an entirely new set of evaluation criteria. We'll then review the current market leaders in cloud based data warehousing to see how they compare.

The options reviewed include:-

- AWS - Amazon Redshift
- Hadoop based Data Warehouse
- Microsoft Azure SQL Data Warehouse
- Oracle Exadata Cloud Service
- Snowflake Elastic Data Warehouse



3. The Benefits of Cloud based Analytics

The following benefits are found in most cloud based analytic solutions:-

- **Infinite Scale:** Unlike on-premises based solutions which are limited by physical data centre rack space, a cloud based solution can (in theory), scale to an almost infinite size. In reality, as each solution has practical limits, these will be indicated where applicable.
- **Low Entry Point:** As cloud based deployments are paid on a subscription or pay-as-you-use basis there's no large capital expenditure to finance, a significant advantage for small to medium sized enterprises.
- **Always Available:** Every solution reviewed includes built in high availability. This compares well to on-premise solutions which typically need an off-site disaster recovery data centre with corresponding additional capital expense.
- **Cost Control:** The pay-as-you-go model supports the ability to control costs on an on-going basis, and avoids a significant capital expense.
- **Time to Market:** Most cloud based databases can be deployed to new hardware within minutes, which speeds delivery of new solutions.
- **Agility:** Many cloud based solutions can be quickly deployed, scaled up or temporarily suspended which makes them an excellent solution for prototyping, proof of concept and discovery analytics.
- **Reduced Cost:** The ability to quickly deploy or remove a database helps reduce costs by avoiding permanently allocated performance and user acceptance test databases.



3.1 Evaluation Criteria

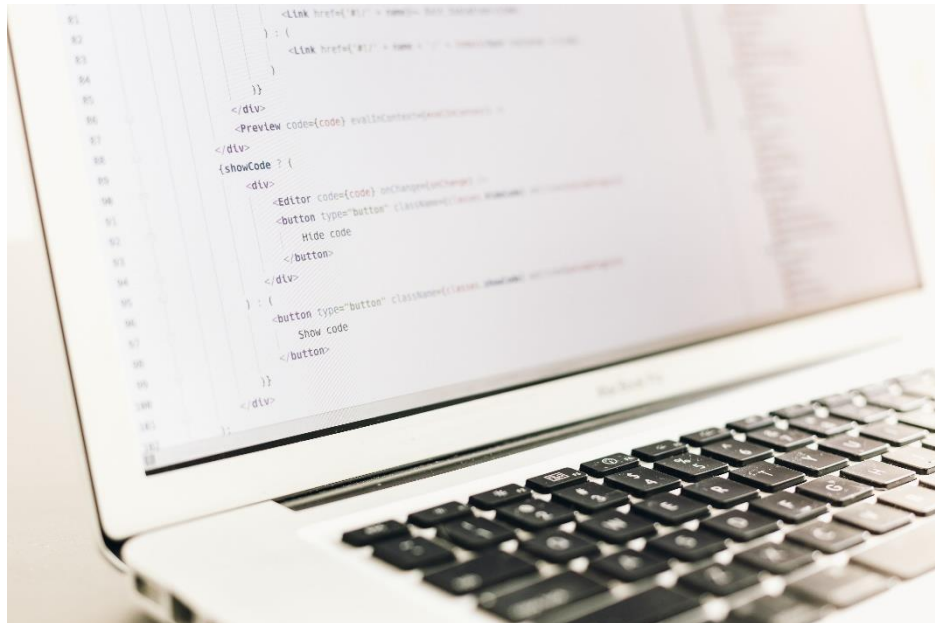
3.1.1 Overview

The relational database has been established for nearly 40 years, and is typically evaluated on performance and usability. However, recent innovations in both database architecture and cloud based deployments, open up an entirely new set of criteria against which to judge a solution.

These include:-

- **Simplicity and Reduced Skills Required:** Modern cloud based solutions are often remarkably easy to deploy and manage, especially when compared to the traditional on-premises solution.
- **On the Fly Elasticity:** Which includes the ability to scale up or down the compute resources to match demands. Some solutions are more quick to respond and can be automatically scaled.
- **Concurrency: Ability to Segment Warehouse Load:** The ability to avoid the *tug of war* between batch ETL processing, and low latency end user queries. This supports a high level of concurrency.
- **Independently Scale Compute and Storage:** Some cloud based solutions can expand storage independent of compute resources. This provides a greater degree of flexibility in deployment, and means you only pay for the resources you need.
- **Management of Diverse Data:** Including the ability to load and query semi-structured JSON data often found in web based applications.
- **Advanced Column Based Storage:** Which has a massive impact upon analytic query performance while maximising data compression.
- **Strong SQL Support:** Which includes analytic window functions and User Defined Functions (UDFs).

The remainder of this report will describe these criteria in detail, and review how the current market leaders compare.



3.2.1 Simplicity and Reduced Skills required

Ten years ago Dr. Michael Stonebraker released the research papers [One Size \[No Longer\] Fits All](#) and [The End of an Architectural Era](#) in which he argued the legacy database solutions from Oracle, Microsoft and IBM were no longer fit for purpose.

Built upon a row based OLTP architecture devised nearly 40 years ago, these have since been extended to include a wide range of features for data warehousing. He argued, these over-complex, general purpose solutions would be replaced by dedicated systems, each designed for a single purpose.

As an example of the complexity, Oracle11g includes a huge range of performance tuning options added over time, including nearly 500 individual tuning parameters and 16 types of index.

Moving to the cloud opens up a genuine opportunity to simplify or even eliminate the system administration burden on these legacy systems. Not only the hardware and operating system maintenance, but the database itself.

To this end, some analytic solutions have been redesigned from scratch, and deliver an innovative architecture which almost entirely removes the need for highly skilled technical resources. This means valuable database engineers and architects are freed up from the demands of database administration to focus on the real challenge, delivering customer insights.

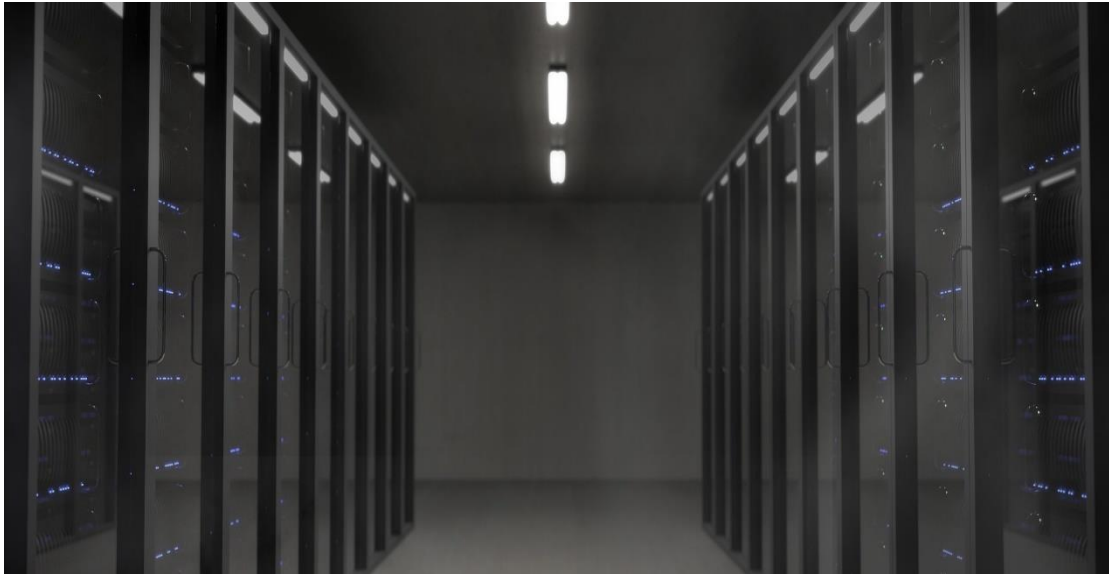
Examples include:

- **Index Design:** To maximise query performance of look-ups. This can be a challenging area as the designer must balance the often conflicting performance requirements of read queries, with the need to rapidly load large data volumes. This often leads to disabling indexes during large batch loads, only to find online query performance suffers.
- **Statistics Data Capture:** To support a cost based optimiser, most databases include tools to analyse and capture metadata statistics to maximise query performance. However, as additional entries are inserted, these statistics can become out of date, which leads to a housekeeping strategy to regularly refresh statistics which can be a CPU demanding task.

- **Horizontal Partitioning and Replication:** Which involves specifying a sensible data distribution method to shard (spread), the data across multiple servers. Although still used by some cloud based Massively Parallel Computing (MPP) systems, this has in some cases been replaced by an entirely new innovative architecture which removes this additional design effort.

Many of these design tasks are eliminated by cloud based solutions. Some also deliver additional features including:-

- **Zero copy cloning:** A technique used to quickly replicate a database to build a fully populated test environment. This works by cloning the database without any physical data copy. This can ease the burden of DEVOPs, as terabytes of data can be cloned within seconds with subsequent inserts and updates allowed on the new data-set.
- **Data Sharing:** Which provides access to both compute and data resources to external partners or subsidiaries on a read-only basis. This avoids the need to build multiple Extract Transform and Load (ETL) pipelines to external users, and avoids the need for Change Data Capture (CDC) routines when the warehouse data is updated, as users always view the latest data.



3.3.1 On-the-Fly Elasticity

On-premise solutions provide a fixed hardware platform with few options for elasticity – the ability to add or remove compute power and storage as needed. An ideal cloud based solution would provide a flexible hardware platform which grows and shrinks as the workload and user numbers change over time.

However, not all cloud based solutions have the same features which should include:-

- **Immediate Response:** The ability to quickly allocate and remove compute resources as needed without any interruption to service, or impact upon performance
- **Scale Up or Down:** The ability to automatically add or remove hardware resources as demands change
- **Small Increments:** The ability to add additional resources in relatively small increments instead of huge step-changes
- **Scale Independently:** The ability to add (or remove) storage and compute resources independently to match resources to current demands.



3.4.1 Concurrency: Ability to Segment Warehouse Load

On premise systems tend to be a fixed size with machine resources shared by disparate groups of users, often with massively varying processing requirements.

Batch ETL processes for example, tend to include long running bulk insert or update operations, and have a very different compute profile to dashboard users who require a fast response time for online analytic queries.

This feature refers to the ability to maximise concurrency by partitioning or segmenting processing to avoid the *tug of war* for machine resources. Ideally, the solution should allow different user groups or workloads to run independently and in isolation while sharing the same data.

Some cloud based solutions however take the traditional approach of placing users a resource queue. However, this leads to a less efficient use of machine resources as capacity is shared between multiple competing user groups.



3.5.1 Independently Scale Compute and Storage

The standard hardware architectures since the 1980s can be categorised by their ability to scale. These include:-

- **Shared Memory:** Describes a system in which everything (including memory), is shared. The option available to scale the solution is to purchase and migrate to a larger platform, and these solutions need a potentially expensive standby machine to support high availability.
- **Shared Disk:** Describes a system in which compute resources are closely connected into a cluster, but share a single the storage resources, typically implemented using a Storage Area Network (SAN).
- **Shared Nothing:** Also known as Massively Parallel Processing (MPP), distributes both processing and data over a number of machines. Data is allocated to each node using a round-robin or consistent hashing method, and additional nodes can be added over time. Although Teradata and Netezza successfully produced data warehouse appliances on this architecture, it does have drawbacks in that processing and storage are tied, and it's not possible to scale them independently.

Recent database research has produced some innovative architectures which challenge the perceived wisdom of the past 40 years. In particular, providing the ability to scale compute processing and storage resources independently.

This leads to the following benefits:-

- **Flexible Scaling:** It is important to be able to flexibly add storage and compute resources as demands change over time. This makes it possible to scale from gigabytes to terabytes and petabytes while independently adjusting the compute resources. This means you can support complex, low latency resource intensive queries on terabyte sized systems, or massive batch processing operations on petabyte sized storage on the same solution.
- **Pay only for what you need:** An ideal solution would allow you to scale storage independent of compute resources. For example, MPP systems are in theory balanced in that adding machines adds both storage and compute resources. However, because these are so closely tied, if storage demands exceed the need for processing capacity, the cost per terabyte can rise disproportionately.

Separating storage from compute resources, means you can scale either independently which means you only pay for the storage and appropriate computing resources you actually need.

“New cloud architectures and infrastructures challenge conventional thinking about big data and colocated storage and compute”.

[Tripp Smith](#) – CTO Clarity Insights.



3.6.1 Manage Diverse Data

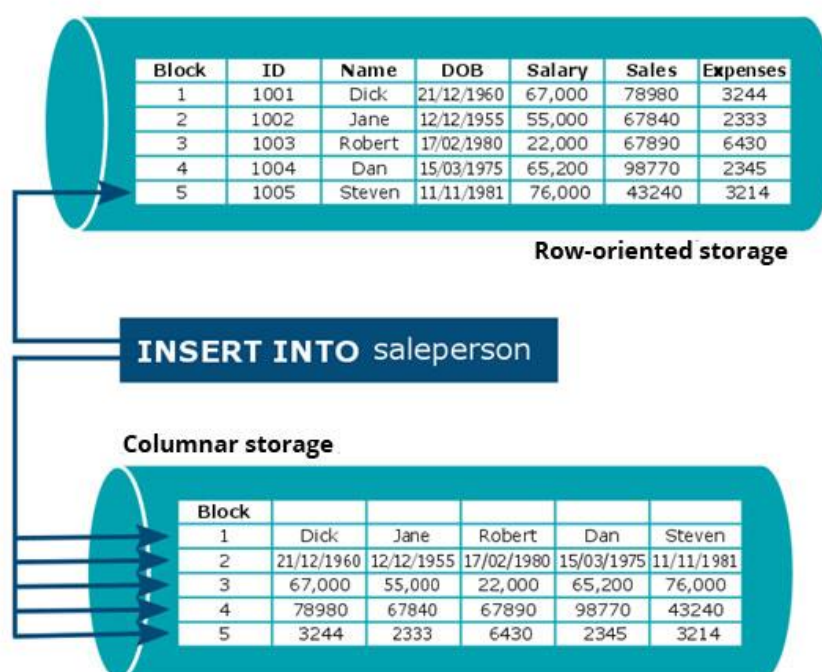
One of the advantages touted by Hadoop HDFS is the ability to handle diverse data including structured data from database sources, and semi-structured data including JSON and XML.

However, it's increasingly possible to parse, store and process JSON and XML directly in the warehouse, and this is especially important when analysing social media feeds and machine sensors or data transferred from web applications.

Some cloud based solutions include SQL extensions to support joining semi-structured JSON data to existing structured data.

Ideally a modern warehouse solution should capture, store and query data against semi-structured data without the need to first parse and translate the data to a predefined relational structure. Given the rise in demand for a Data Lake solution, there are also benefits in the ability to store data in (as yet) undetermined data format to support schema on read.

3.7.1 Advanced Column based Storage



The traditional databases from Oracle, Microsoft and IBM extend an existing *row-based* solution designed for OLTP applications with additional features for a data warehouse. However, [extensive research](#) has since demonstrated that column based databases are as much as 100 times faster for analytic queries, and most (but not all) modern data warehouse solutions are built upon a columnar data store.

The diagram above illustrates the difference between the two solutions. The row based storage above works well for OLTP systems which typically update all columns on a single row. However, analytic queries tend to summarise or group by a limited number of columns over millions of rows.

In the example above, a sum of EXPENSES by NAME would query just two columns, and on a wide fact table, the performance gains can be significant.

Storing data in columnar format also means there's greater opportunities to compress the data by removing repeating values, which can lead to significant cost savings.

"In most real-world environments, column stores are 50-100 times faster than row stores"

- [Dr. Michael Stonebraker \(MIT\)](#)



4. Analysis of Options

The cloud analytic warehouse solutions from the following vendors were evaluated (in alphabetic order):-

- Amazon Redshift
- Hadoop based Data Warehouse
- Microsoft Azure SQL Data Warehouse
- Oracle Exadata Cloud Service
- Snowflake Elastic Data Warehouse

Each was evaluated based upon the predefined criteria, and the strengths and weaknesses of each are explored.



4.1 Amazon Redshift

4.1.1 Overview

In common with the Oracle and Microsoft solutions, *Amazon Redshift* was not built specifically for cloud deployment, but was based upon an existing on-premises data warehouse solution, ParAccel, licenced from Actian in 2011.

Redshift offers a massively parallel processing (MPP) architecture over a column-oriented database to support fast execution of analytic queries against terabytes to petabytes of data. Customers can deploy a warehouse within minutes on a dense compute or storage nodes starting from as little as 160Gb SSD, rising to 2 petabytes of disk storage on a 128 node cluster.¹

Although not evaluated here, it is worth noting that Amazon recently released an extension to Redshift called *Spectrum* which supports on-the-fly elasticity and separation of compute and storage. However, from an initial viewing, this product appears to lack maturity.

“AWS has been an excellent technology partner for all our data needs”
[Gartner Peer Insights Review. 2017.](#)

4.2.1 Simplicity and Reduced Skill-set required

Although Redshift is a hosted data warehouse platform, a [survey published in January 2017](#) at **Amazon re:Invent** indicated that almost 60% of respondents indicated the data warehouse was difficult to manage, although 51% were satisfied with the performance.

However, an analysis of over 70 customer reviews on [Gartner Peer Insights](#) indicates a high level of overall satisfaction with the solution, and 80% of customers would recommend AWS to other clients.

It's worth pointing out, Redshift is a pure MPP solution in which data is distributed across multiple nodes in a cluster using a *distribution key*². It is therefore important to carefully design the schema to maximise data co-location and query performance. It's not necessary to manage indexes which adds to the simplicity, although statistics gathering is a separate manual process.

Finally, Redshift is a fully managed data warehouse service with automatic data backups, upgrades and patches applied. However, although there are compelling advantages, Redshift will still need technical management and some DBA support.

¹ [Amazon Redshift Clusters](#)

² [Choosing a Distribution Key](#)

4.3.1 On-the-Fly Elasticity

Although Redshift does not support on-the-fly elasticity, it is possible to resize³ a cluster. However, this requires a cluster reboot which places the source machines in read-only mode while copying data to the new cluster. The read-only database remains open for queries, although this may extend the time taken to migrate the data.

The alternative method⁴ takes a database snapshot to Amazon S3 storage, and restores the data to a target cluster, where it can be resized before operations are switched over. This has the advantage the cluster remains open for write operations during the resizing operation.

All resizing operations are a manual step, and there's currently no automatic way to adjust machine resources to match workload.

4.4.1 Ability to segment the warehouse load

Running mixed workloads is managed by the Workload Manager (WML) which defines a series of user groups designed to segment workloads by priority. Each group is assigned a different workload queue⁵ to ensure fast running interactive queries don't get stuck in queues behind long running ETL tasks.

In addition, each queue can be assigned up to eight WML rules⁶ to prevent long running queries on a low latency queue. This can be used for example, to abort queries on a short queue after a predefined elapsed time.

4.5.1 Independently Scale Compute and Storage

As Redshift is a scale out MPP solution, it is not possible to independently scale storage and processing, as each node adds both compute and disk based capacity.

4.6.1 Manage Diverse Data

Redshift provides a limited number of SQL functions to manage semi-structured JSON data. Functions include the ability to validate a JSON document or extract elements, but only VARCHAR data type is supported meaning there's no support for documents larger than 64K which may be a severe limitation in some cases.

³ [Resizing a cluster](#)

⁴ [Snapshot, Restore and Resize a cluster](#)

⁵ [Workload Management](#)

⁶ [WLM Query Rules](#)

4.7.1 Advanced Column based Storage

One of the greatest performance challenges of an analytic data warehouse is reducing the disk accesses and data transfer times. Redshift maximises query performance by storing data in columnar format. As most analytic queries access a few columns from relatively wide tables, this can lead to dramatic performance improvements over legacy row based storage.

Redshift stores data in columnar format in 1MB blocks, and records the min and max value of each column. Combined with the use of *sort keys*, this means Redshift can further improve performance by filtering blocks depending upon the SQL where clause.

4.8.1 Observations

Unlike some warehouse solutions, it is not possible to suspend a cluster once started, and billing continues even when the cluster is not being used. It is possible to stop a service, but the database must be backed up and deleted before later being reloaded from S3 storage.

Redshift provides a migration path for customers using Postgres as the underlying technology is built upon PostgreSQL version 8. This is however based upon the 2005 release, and the full range of operations is not available including stored procedures and referential constraints.

Finally, although Redshift appears easier to manage than Oracle or Microsoft Azure, it wasn't completely architected for the cloud, and lacks the ease of management of Snowflake.

*“We imported 20 million rows of sales facts in less than 15 seconds,
and the content was query-able immediately”*

- [Customer Review - IT Central Station](#)



4.2 Hadoop based Data Warehouse

Although not by any standards *a database*, it's useful to consider the limitations of Hadoop, if only to rule it out as a contender as an analytic warehouse.

Originally developed by Doug Cutting based upon the [Google File System \(GFS\)](#) and [MapReduce](#) papers published in 2003-4, this was an attempt to deliver a resilient but simple batch processing solution to run on a scale out architecture using commodity servers, and directly attached storage.

Despite the addition of interactive SQL query engines including Cloudera Impala, Presto and Hive, the solution has not managed to achieve the performance expected of a database, and the query optimisers remain immature.

Finally, HDFS allows no control for physical data placement and co-location found on MPP systems, and our experience with Hadoop has demonstrated, while it's possible to achieve good batch query performance on large data volumes using brute force, the solution is hard to scale and makes inefficient use of resources, especially as additional users are added.

“Most of those who expected Hadoop to replace their enterprise data warehouse have been greatly disappointed” - [James Serra \(Microsoft\)](#).

4.1.2 Simplicity and Reduced Skill-set required

Unlike the alternatives proposed, Hadoop is not a single product or even a collection of products from a single vendor, but an ecosystem of open source software. This means deploying a Hadoop based solution requires knowledge and skills in a potentially huge number of tools including HDFS, Spark, Impala, Hive, Yarn, Flume, Sqoop, Zookeeper and Kafka.

Hadoop installation and administration is a complex process, and it is a challenging solution to architect and deploy. It's probably only a sensible option for an experienced IT department with extensive skills in Java and distributed systems architecture.

4.2.2 On-the-Fly Elasticity

Although Hadoop provides theoretically infinite scalability that only extends to adding nodes to the cluster, and there is no support of on-the-fly elasticity, even in the cloud.

4.3.2 Ability to segment the warehouse load

It is not possible to physically segment workloads on a Hadoop cluster. Instead workload management uses a scheduler such as YARN or Mesosphere to separate batch ETL and online reports on different queues. However, resource allocation is basic and mainly based upon data volumes processed.

4.4.2 Independently Scale Compute and Storage

The Hadoop architecture relies for its performance upon directly attached storage on a cluster of machines, and although it is [possible](#) to run Hadoop on independent SAN storage, it's not advisable. This means you need to scale the compute resources in line with storage.

4.5.2 Manage Diverse Data

As the underlying storage system Hadoop Distributed File System (HDFS) is effectively a file system, Hadoop has strong support for diverse data and HDFS can be used to store structured, semi-structured and entirely unstructured data. This can include sound, video, image and JSON data, and there's a wide range of open source tools (eg. [Apache Solr](#)) to index, search and process diverse data.

It could however be argued that the majority of analytic data in the warehouse is structured, and query performance and scalability over structured data is not a Hadoop strength.

4.6.2 Advanced Column based Storage

Hadoop does support columnar based storage using Parquet and ORC formats, and both support column projection and storage indexes which provide block level partition elimination to maximise performance. Unlike Redshift and Snowflake however, they don't support sort keys which would greatly improve performance.

4.7.2 Observations

Although originally promoted as a potential replacement for large expensive data warehouse appliances, Hadoop has proven to fall short of the performance and scalability required, especially for business intelligence applications.

As a compute platform, it is remarkably complex to deploy and administer, and although a potential platform for an on premises Data Lake, it's not a suitable replacement for a modern cloud based analytics database.



4.3 Microsoft Azure SQL Data Warehouse

4.1.3 Overview

Based upon the on-premises *Microsoft Parallel Data Warehouse* platform, this solution blends a massively parallel processing (MPP) architecture as used by Redshift, with the separation of compute and storage seen in Snowflake.

However, unlike both these alternatives this solution includes a great deal of inherited complexity, and like Oracle Exadata will require a team of skilled Database Administrators and database engineers to operate and tune. It's also a relatively new product, and does not, for example include a sophisticated SQL query optimiser found in the other solutions.

“Microsoft has made it as painless as possible to implement a data warehouse, especially when compared to deploying a solution on-premises”

– [Robert Sheldon](#) (Redgate Hub).

4.2.3 Simplicity and Reduced Skill-set required

Microsoft Azure Data Warehouse has much in common with Oracle Exadata in the cloud, in that they were both based upon legacy on-premises database solutions. This clearly shows in the Azure documentation which is extensive, and includes a long and complex cheat sheet⁷ including detailed advice on how to manage data distribution and replication, and strategies for indexing, partitioning and statistics maintenance.

As such, this solution has none of the advantages of simplicity or reduced system management cost associated with some of the alternatives.

4.3.3 On-the-Fly Elasticity

It's relatively quick and easy to scale the system as compute and storage services are separated, and there's even an option to pause compute resources which helps control costs.

However, while it's relatively easy to pause or resume the warehouse using the web based Portal, this does not have the full in flight control of Snowflake, and any currently executing read queries are immediately terminated.

⁷ [SQL Data Warehouse Cheat Sheet](#)

4.4.3 Ability to segment the warehouse load

In common with Oracle, Microsoft uses a workload management⁸ approach to segment and control load. Using this method, the DBA defines a number of *resource classes* which can be assigned a database role with associated service level defined in Data Warehouse Units.

Using this method, a database administrator can control compute resources allocated to different user groups each with a different workload profile, although in practice this is likely to be an on-going challenge, resolved by trial and error. Also once the concurrency limit is reached, queries are queued.

4.5.3 Independently Scale Compute and Storage

Microsoft does provide separation of compute and storage resources, with data stored and backed up from Azure Blob storage, and customers are charged separately for storage and compute resources.

Compute resources are provided by compute nodes, and unlike the alternatives which allocate a fixed number of machines, Microsoft opts for an abstract billing method called Data Warehouse Units (DWUs) which represent an aggregate bundle of CPU, I/O and memory resources.

For higher performance and increased concurrency (more queries executed in parallel), the user increases the number of DWUs which are billed separately from storage.

4.6.3 Manage Diverse Data

In common with other options, Azure Data Warehouse provides SQL extensions⁹ and functions to manage data in JSON format, and this can be used to combine structured and semi-structured data. These extensions support outputting relational data in JSON format, and querying or modifying JSON in the database.

4.7.3 SQL Language

This platform uses Transact-SQL, and is accessible from familiar tools including SQL Server Management Studio (SSMS). It is however, not a complete port, and there are extensive differences¹⁰ and unsupported¹¹ features.

In conclusion, while TSQL provides a head-start, it's possible there will be challenges around migrating an existing data warehouse to the Azure cloud platform. Customers should also be mindful of performance challenges until the issues around query optimisation and stats gathering are resolved.

⁸ [Microsoft Azure workload management](#)

⁹ [Managing JSON in Azure SQL database](#)

¹⁰ [Differences from SQL Database](#)

¹¹ [Common T-SQL Limitations](#)



4.4 Oracle Database Cloud Exadata Service

4.1.4 Overview

First released as an on-premises database machine, Exadata provides a combined compute and storage server connected by a fast InfiniBand network connection, and was introduced in the cloud in 2015 as a managed service.

Exadata is provided as part of a range of offerings including outright purchase or on a subscription basis on premise. This review describes the cloud based managed service.

Oracle have announced the intention to deliver a [fully autonomous](#) data warehouse cloud solution, although as yet, few details are available.

4.2.4 Simplicity and Reduced Skill-set required

As the Exadata Cloud Service is effectively a database as a managed service, it only eliminates the effort around the deployment and management of hardware. Unlike the solutions provided by Snowflake and Redshift, Exadata requires a DBA to manage databases, and a team of designers to complete physical database design, performance tuning and system monitoring.

4.3.4 On-the-Fly Elasticity

Although it's possible to expand an existing cluster by adding more compute resources, it's a manual operation involving activating additional cores. Oracle do however provide a "Compute Bursting" option up to double the processor capacity charged at a metered hourly rate to allow for peak usage.

This provides some flexibility over an on-premises Exadata solution which would normally be sized for the highest possible workload, although not as flexible as the other cloud based solutions available.

4.4.4 Ability to segment the warehouse load

Although Oracle deploys a minimum quarter rack which means users run on their own dedicated hardware, there are few options available to quickly adapt the solution to segment workloads.

4.5.4 Independently Scale Compute and Storage

Although there's some options to scale the storage and compute servers, these are tightly coupled, and the quarter, half or full rack servers come with a fixed range of 3, 6 or 12 storage servers ranging from 42-168TB on the X5-2 service.

4.6.4 Manage Diverse Data

Oracle has supported CLOB data types, (character large object), for some time, and now includes support for semi-structured data including the ability to insert, validate and query JSON data. Oracle provides a number of functions for use in SQL to extract JSON elements in addition to loading JSON files as external tables, and indexing documents.

4.7.4 Advanced Column based Storage

Although Oracle describe the Exadata solution as using Hybrid Columnar Compression, it is not a full columnar storage system, although it does gain some performance benefits from column projection on the storage server.

Oracle does however store data in columnar format in memory if the Database In-Memory option is used. This allows the DBA to select a predefined sub-set of tables, columns and partitions, which will be locked in memory in a separate columnar cache, complete with storage indexes to maximise analytic query performance.

4.8.4 Observations

Oracle Exadata Cloud service is effectively a *managed service* version of the on-premises hardware. It loses out when compared to other cloud based options, and provides none of the agility, elasticity and infinite scalability of the alternatives including Snowflake.

Oracle does however benefit from a large installed base of corporate customers with a significant investment in application code and skills, and the Exadata solution is 100% code compatible with existing Oracle systems.

For customers heavily invested in Oracle, this makes the migration path to the cloud a much easier step than switching to another database, and the cloud based service comes with the advantage that all Oracle options are included, whereas on-premise they're charged on an individual option basis.

It's likely however that many customers will take the opportunity to re-evaluate the technology stack, and consider the alternatives, especially if this helps avoid vendor lock in, and potentially reduce costs.

Customers should be aware, that unlike Snowflake that will run on multiple cloud vendors, Oracle like Microsoft and Amazon Redshift is only available from a single cloud based platform.

*“[A legacy row based] RDBMS represents a major potential bottleneck
in the data warehouse”*

— [Mark Lewis. Clarity Insights.](#)



4.5 Snowflake Elastic Data Warehouse

4.1.5 Overview

Founded in 2012, Snowflake was designed by a team headed by ex-Oracle architects to deliver a cloud analytics solution from scratch. Unlike some other vendors who migrated on-premise based solutions, Snowflake was able to introduce some quite innovative architectural solutions, and exploit the potential benefits of a cloud based solution.

Deployed as a multi-cluster solution with complete separation of compute and storage, Snowflake provides a solution which runs completely independent workloads against a shared data resource without contention, and supports a highly flexible elastic compute resource which can be scaled up or down within seconds.

“With superior performance and the most hands-off model of ownership, Snowflake is the epitome of a data warehouse as a service”

— [William McKnight. Gigaom Research](#)

4.2.5 Simplicity and Reduced Skill-set required

Snowflake includes an innovative architecture approach which reduces complexity and management overhead. Unlike typical MPP solutions, there’s no need to shard data across nodes, and adjusting the compute resource is an almost instant operation that on the Enterprise version can be automatically triggered¹² to match the workload.

Similar to Redshift, Snowflake has no indexes to create or maintain, and no horizontal or vertical data partitioning to specify, and data is stored in micro-partitions which makes use of partition elimination to maximise performance.

All queries are optimised prior to execution, and there are no data statistics to maintain. Database backup, and cross-datacentre disaster recovery is automatic, as is the ability to re-query and restore data modified up to 90 days ago.

Finally, the solution includes the ability to produce instant zero copy clones of entire databases which aids agility, and data may be securely shared with external partners.

In short, Snowflake is built for ease of use, and all but eliminates the need for highly skilled database administrators, and data engineers and architects are freed up to work on logical design and development rather than tuning the solution for performance.

¹² [Snowflake Automatic concurrency scaling](#)

4.3.5 On-the-Fly Elasticity

Snowflake supports near instant elasticity both to a larger or smaller cluster, and it's possible to suspend a cluster when not in use. As a result of the independent design of compute and storage nodes, processing can be migrated to a larger or smaller cluster instantly, with all new queries starting on the new machine.

Finally, the ability to suspend a cluster when not in use means there's more options to control costs which is further enhanced by per-second billing.

4.4.5 Ability to segment the warehouse load

Unlike some solutions in which workload is managed by different queues, Snowflake allows multiple independent clusters of different sizes to access the same storage pool. This means a cluster can be reserved for large batch ETL processing, and run independently of low-latency dashboard queries or compute intensive data science prototypes.

This effectively eliminates contention for machine resources as each user group can have their own *virtual warehouse* running on separate hardware.

4.5.5 Independently Scale Compute and Storage

Snowflake is deployed as a series of one or more clusters of varying sizes, each backed by a set of locally attached solid state disks for caching. These compute resources are completely independent of each other, and access a common storage resource on Amazon S3, which means it's possible to independently size storage and compute resources.

4.6.5 Manage Diverse Data

One of the main drivers behind NoSQL databases like MongoDB and Couchbase was their ability to natively handle semi-structured data including JSON and XML. This does however have the disadvantage of adding yet another database to the toolset, and requires additional skills to manage.

Snowflake includes support for a new *variant* data type to natively store JSON data while automatically optimising it in the background for columnar access. It also supports SQL extensions to query and flatten JSON data, which allows users to combine structured and semi-structured data within the same query.

4.7.5 Advanced Column based Storage

Snowflake stores data in column based format to maximise query performance and data compression. This means only the columns needed are read from disk, and data is processed in compressed form which avoids the need to decompress data upon read.

“Combining Snowflake’s unique cloud data warehouse with Looker, business teams and data analysts can curate and explore data at any scale with outstanding performance and flexibility.”

– [Frank Bien. CEO, Looker](#)

5. About Sonra Intelligence

Sonra Intelligence are experts in data warehouse design, implementation, and cloud migration. We provide services across the globe including [Redshift](#) and [Snowflake](#) consulting, [training](#) and [data architecture advisory services](#).

With offices in Dublin (Ireland) and Stettin (Poland), we provide advice and guidance to blue chip clients on a range of technologies including Amazon Redshift, Snowflake, Oracle, and Hadoop.

Sonra is also developer of Flexter. [Flexter](#) is an ETL tool for XML and JSON. It automatically converts XML/JSON to text, a relational database, or Hadoop. Visit the [FAQ](#) section for common queries.

[Reach out to us](#), and see how we can help you modernise your data analytics platform.

Web: <https://www.sonra.io>
eMail: hello@sonra.io
Telephone: +353 1 5345 015

6. About the author: John Ryan



John Ryan is an experienced Big Data and Data Warehouse Architect with over 30 years IT and consultancy experience in a range of industries as diverse as Mobile Telephony, Pharmaceuticals and Financial Services.

Author of several eBooks on database technology and industry trends, he regularly writes articles for platforms including ODBMS Magazine, Medium and LinkedIn. You can follow John on his blog [Analytics Today](#).

7. About the co-author: Uli Bethke



Uli is a thought leader in the data industry and has architected and delivered data projects in Europe, North America, and South East Asia. He is a regular speaker at conferences, holds an Oracle ACE award, and has written and reviewed various books.

Uli is the founder of the [Hadoop User Group Ireland](#). He is also a co-founder and VP of the Irish chapter of DAMA, a non for profit global data management organization.

8. Training Big Data for Data Warehouse Professionals

Fix common data warehouse headaches with big data concepts and technologies.
Find out about the limitations of big data.

Duration: 1 day. Onsite in EMEA. Virtual classroom webinar world-wide.

Trainer: Uli Bethke, CEO Sonra.

These companies are all Big Data certified



What's the trouble with the data warehouse?

- Data warehouses are bursting at the seams.
- Data volumes and costs are exploding.
- Relational data warehouses don't get on with graphs, unstructured data, keyword searches, machine learning, predictive analytics etc.
- It takes forever to get data into the data warehouse and insights out of it.
- Only 20% of enterprise data finds its way to the data warehouse.
- Writing ETL code is slow and hard to maintain.
- Traditional data warehouse architecture sucks for real-time analytics.

Does this sound familiar?

Yes? Then this training course is for you.

[Find out more](#)