

VERTICA

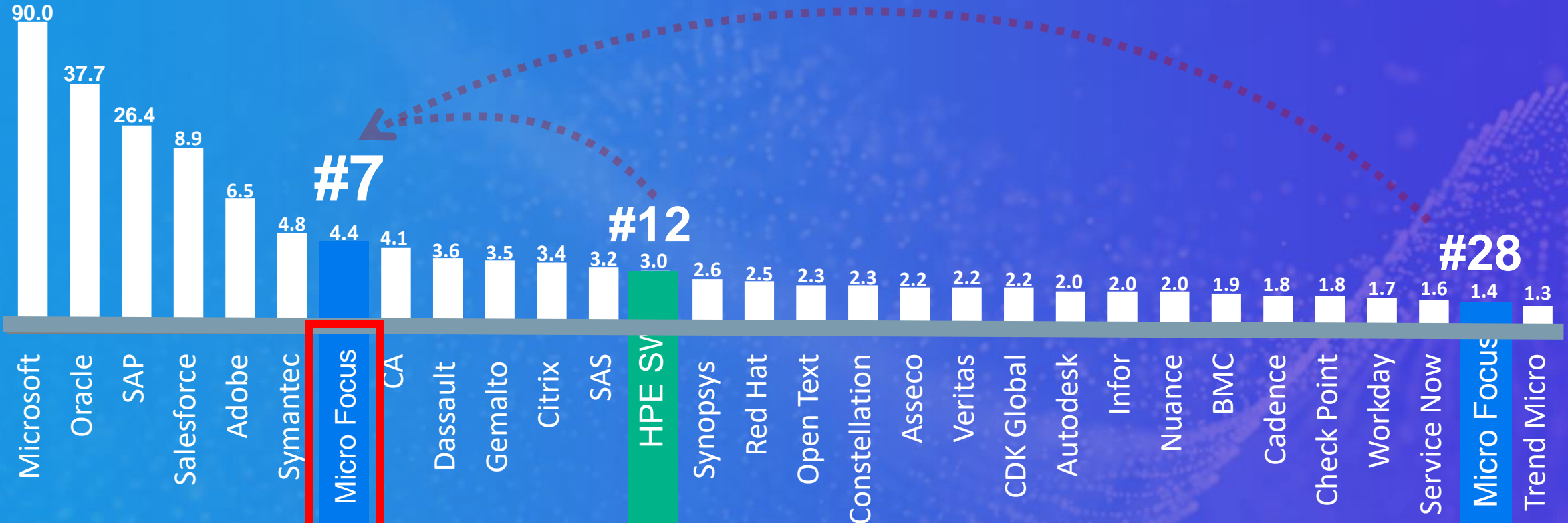
## Vertica Analytics Platform

데이터 기반 비즈니스를 위한 데이터 분석 플랫폼

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

# 글로벌 7위의 SW회사 Micro Focus

2017년 9월 Micro Focus는 업계12위인 SW부분과 업계28위인 Micro Focus가 합병하여 업계 7위의 SW 회사가 되었습니다



# 글로벌 7위의 SW회사 Micro Focus



Hewlett Packard  
Enterprise



COBOL



40  
Years

Network  
Management  
Data Protector



30  
Years



2017.9.1



# 글로벌 최고의 회사들이 분석업무를 위해 Vertica를 사용합니다



**“Digital Darwinism is unkind to those who wait.”**

- Ray Wang, Constellation Research, June 2015



# VERTICA

Turing Award 수상자인 Michael Stonebraker 박사의 C-Store 프로젝트의 결과로 2006년 시장에 출시되었습니다  
Michael Stonebraker 박사는 Greenplum, Netezza 등의 기초인 Postgres DB를 개발한 세계적인 석학입니다.

## C-Store: A Column-oriented DBMS

Mike Stonebraker\*, Daniel J. Abadi\*, Adam Batkin\*, Xuedong Chen†, Mitch Cherniack\*,  
Miguel Ferreira\*, Edmond Lau\*, Amerson Lin\*, Sam Madden\*, Elizabeth O'Neil\*,  
Pat O'Neil†, Alex Rasin‡, Nga Tran\*, Stan Zdonik‡

\*MIT CSAIL  
Cambridge, MA

\*Brandeis University  
Waltham, MA

†UMass Boston  
Boston, MA

‡Brown University  
Providence, RI

### Abstract

This paper presents the design of a read-optimized relational DBMS that contrasts sharply with most current systems, which are write-optimized. Among the many differences in its design are: storage of data by column rather than by row, careful coding and packing of objects into storage including main memory during query processing, storing an overlapping collection of column-oriented projections, rather than the current fare of tables and indexes, a non-traditional implementation of transactions which includes high availability and snapshot isolation for read-only transactions, and the extensive use of bitmap indexes to complement B-tree structures.

We present preliminary performance data on a subset of TPC-H and show that the system we are building, C-Store, is substantially faster than popular commercial products. Hence, the architecture looks very encouraging.

### 1. Introduction

Most major DBMS vendors implement record-oriented storage systems, where the attributes of a record (or tuple) are placed contiguously in storage. With this *row store* architecture, a single disk write suffices to push all of the fields of a single record out to disk. Hence, high performance writes are achieved, and we call a DBMS with a row store architecture a *write-optimized* system. These are especially effective on OLTP-style applications. In contrast, systems oriented toward ad-hoc querying

in which periodically a bulk load of new data is performed, followed by a relatively long period of ad-hoc queries. Other read-mostly applications include customer relationship management (CRM) systems, electronic library card catalogs, and other ad-hoc inquiry systems. In such environments, a *column store* architecture, in which the values for each single column (or attribute) are stored contiguously, should be more efficient. This efficiency has been demonstrated in the warehouse marketplace by products like Sybase IQ [FREN95, SYBA04], Addamark [ADDA04], and KDB [KDB04]. In this paper, we discuss the design of a column store called C-Store that includes a number of novel features relative to existing systems.

With a column store architecture, a DBMS need only read the values of columns required for processing a given query, and can avoid bringing into memory irrelevant attributes. In warehouse environments where typical queries involve aggregates performed over large numbers of data items, a column store has a sizeable performance advantage. However, there are several other major distinctions that can be drawn between an architecture that is read-optimized and one that is write-optimized.

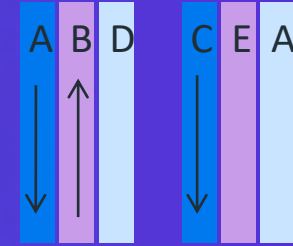
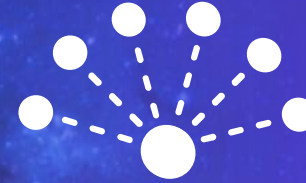
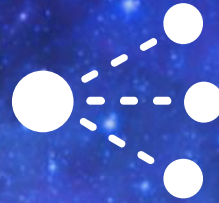
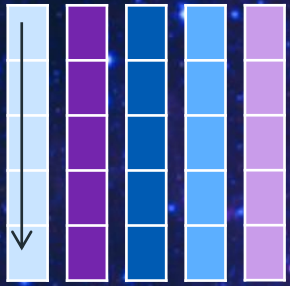
Current relational DBMSs were designed to pad attributes to byte or word boundaries and to store values in their native data format. It was thought that it was too expensive to shift data values onto byte or word boundaries in main memory for processing. However, CPUs are getting faster at a much greater rate than disk bandwidth is increasing. Hence, it makes sense to trade CPU cycles, which are abundant, for disk bandwidth, which is not. This tradeoff appears especially profitable in a read-mostly environment.

Mike  
Stonebraker





# VERTICA 5대 기술 요소



## Native Columnar Storage

필요한  
컬럼만을  
조회하여 빠른  
쿼리 성능 보장

## Compression /Encoding

I/O 비용을  
최소화하는  
동시에 성능을  
가속화

## MPP Scale-out

Name node와  
같은 single point  
of failure를  
제거한 순수  
MPP 아키텍처  
Exabyte 수준의  
확장성 제공

## Distributed Query

특정 노드에  
대응하는  
분산 쿼리  
수행

## Projections

노드 장애 대처와  
쿼리 성능을  
담보하기 위한  
최적화 방안 제공

VERTICA 9.3 GA 2019.10.15



# Vertica 포트폴리오



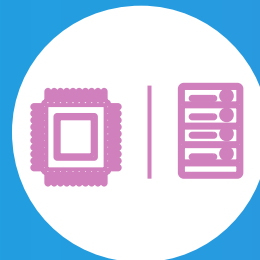
## Vertica Enterprise On-Premise

- 컬럼 처리 및 선진 압축 기법
- 최대 성능 및 확장성
- 다양한 선진 기법 제공  
(Machine Learning, Pattern matching, Flex Tables 등)



## Vertica Enterprise in the Clouds

- 클라우드 플랫폼으로의 빠른 전개
- AWS, Azure, Google, Vmware 지원
- 유연한 클라우드 기반 옵션 제공



## Vertica Eon in the Clouds & On-Premise

- Computing Node 와 Storage 분리
- Object Storage 기반의 무한한 확장성
- Workload에 따른 유연한 증설 및 Multi-Cluster 구성 지원

## 경쟁기술

: Oracle Exadata, Teradata, Greenplum

## 경쟁기술

: Redshift, MEMSQL, Snowflake

## 경쟁기술

: 없음

[Top Rated](#) | [Features](#) | [All Solutions](#)

### Best Data Warehouse Solutions, Comparisons and Vendors

The best Data Warehouse vendors are Oracle Exadata, Teradata, Vertica, Apache Hadoop and Netezza. Oracle is the top solution according to IT Central Station reviews and rankings. One reviewer writes: "We can run the exact same software as the latest x7-2 Exadatas and we can even install virtualization if we want to.", and another reviewer writes: "The reports are always readily available.". The 2nd best product is Teradata. A user writes: "This solution has proven technology with significant advantages in some key areas, and greater depth and experience in providing industry knowledge and solutions.", and another reviewer writes: "Improved the performance of our ETL procedures and reporting". See our free [Buyer's Guide for Data Warehouse](#).

Ranking explanation Company size: All Rankings through: Sep 2019




	1. Oracle Exadata	8.8	12	487	19,215
		Average Rating	Reviews	Words/Review	Views
	2. Teradata	8.2	12	228	19,349
		Average Rating	Reviews	Words/Review	Views
	3. Vertica	8.5	11	272	15,217
		Average Rating	Reviews	Words/Review	Views

[Top Rated](#) | [Features](#) | [All Solutions](#)

### Comparison of Best Cloud Data Warehouse Solutions & Vendors

The best Cloud Data Warehouse vendors are Vertica, Snowflake, Amazon Redshift, Microsoft Azure SQL Data Warehouse and Oracle Autonomous Data Warehouse. Micro Focus is the top solution according to IT Central Station reviews and rankings. One reviewer writes: "Easy to implement, by tuning the model (projection design) you get great performance", and another reviewer writes: "Allows us to take volumes and process them at a very high speed". The 2nd best product is Snowflake. A user writes: "A flexible solution with good clustering, and the pay-per-use feature is useful", and another reviewer writes: "The distributed architecture of Snowflake has the capacity to process huge datasets faster and allows us to scale up and down according to our needs". See our free [Buyer's Guide for Cloud Data Warehouse](#).

Ranking explanation Company size: All Rankings through: Sep 2019

	1. Vertica	8.5	11	272	15,217	8,628
		Average Rating	Reviews	Words/Review	Views	Comparisons
	2. Snowflake	8.7	3	527	15,169	11,128
		Average Rating	Reviews	Words/Review	Views	Comparisons
	3. Amazon Redshift	8.7	3	195	12,646	9,069
		Average Rating	Reviews	Words/Review	Views	Comparisons

# The sun is shining on the Clouds



Google Cloud Platform

Cloud economics offer budget relief for certain workloads

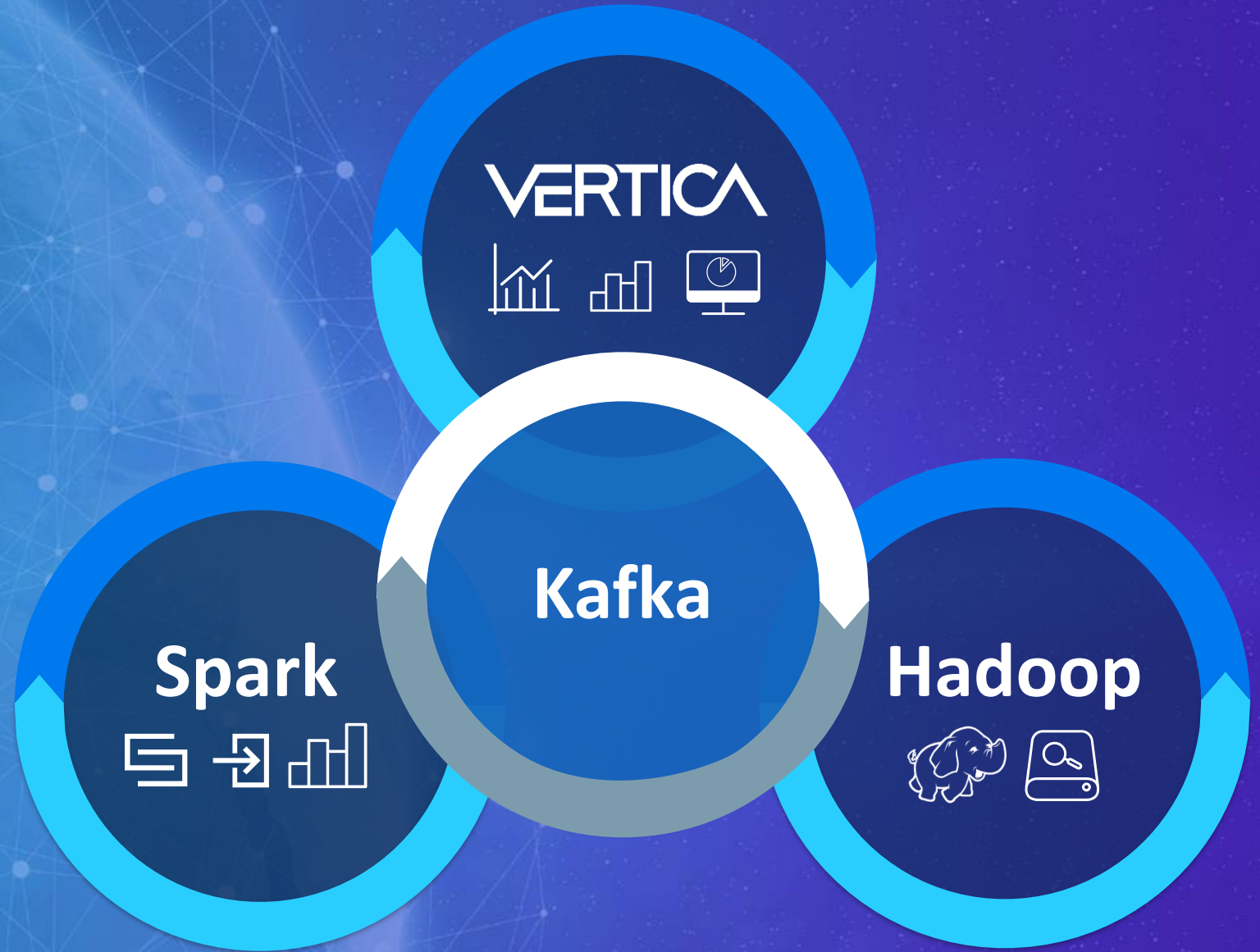
Data storage and compute are the first stop on the journey

Cloud based analytics stacks enter the market



# Embracing an Open Source Architecture

Apache Spark, Hadoop and Kafka Integration



# An Open Architecture Integrated with Rich Ecosystem

## Data Transformation



## Messaging



## ETL



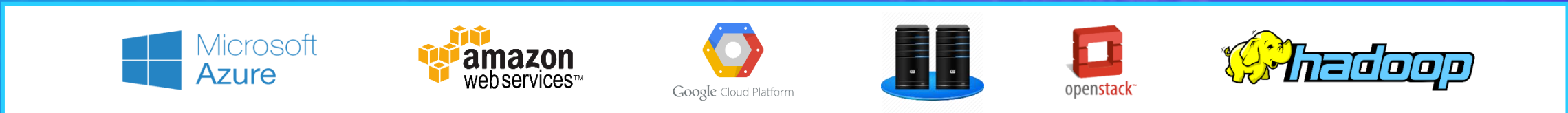
USER  
DEFINED  
LOADS

## User-Defined Functions

R	Java	C++	Python	SQL
Geospatial		Real-Time		Text Analytics
Event Series		<b>VERTICA</b>		Pattern Matching
Time Series		Machine Learning		Regression
User Defined Storage				
Security				
External tables to analyze in place				

ODBC  
JDBC  
OLEDB

## BI & Visualization



# Vertica Roadmap

Vertica는 Micro Focus의 핵심자산으로 지속적인 투자를 바탕으로 한 기술혁신을 주도합니다.

## Foundation

- Columnar Store
- Aggressive Data Compression
- MPP Architecture
- HA Architecture
- ANSI SQL Compliant
- Java, Python, R APIs
- ACID Compliance
- No Single Point of Failure
- Management Console
- Database Designer
- Projections and Optimizations

## Mission



다양한 Cloud 플랫폼 지원



선진화된 In-Database 분석

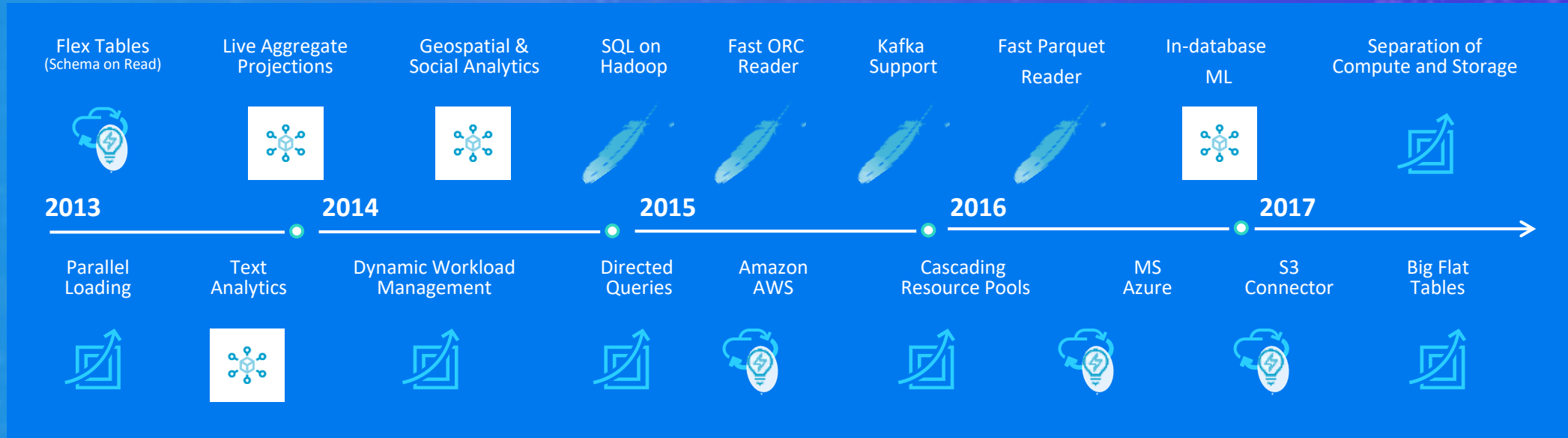


Exabyte 수준의 성능 보장



Open Source 연계

## Innovation Timeline







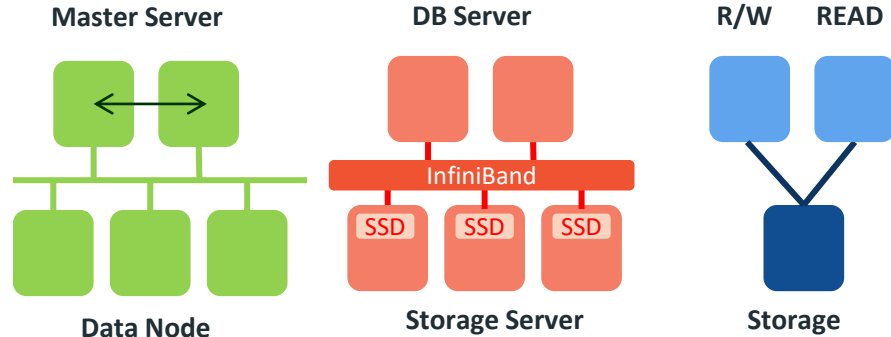
VERTICA

# Technical Advantage

# Vertica 아키텍처의 특징 > Pure-MPP(Massively Parallel Processing)

Vertica는 별도의 마스터 노드를 분리하지 않고 **모든 노드가 동일한 역할을 수행하는 pure-MPP 구조**입니다. 또한 클러스터 구성이나 노드 구성에 특별한 하드웨어나 소프트웨어를 필요로 하지 않기 때문에 비용적인 강점이 있으며 복잡한 구성이 불필요하여 클러스터 구성 과정이 매우 간단하고 빠릅니다.

## Other Systems



### 타사 appliance 시스템

- 두 종류 이상의 서버
- 특별한 H/W 사용으로 복잡도 증가
- 마스터 서버를 통한 작업 수행
- 저가형 서버 사용
- 관리/운영에 다양한 고려사항 존재

## VERTICA



- No specialized nodes
- All nodes are peers
- Query/Load to any node
- Continuous/real-time load & query

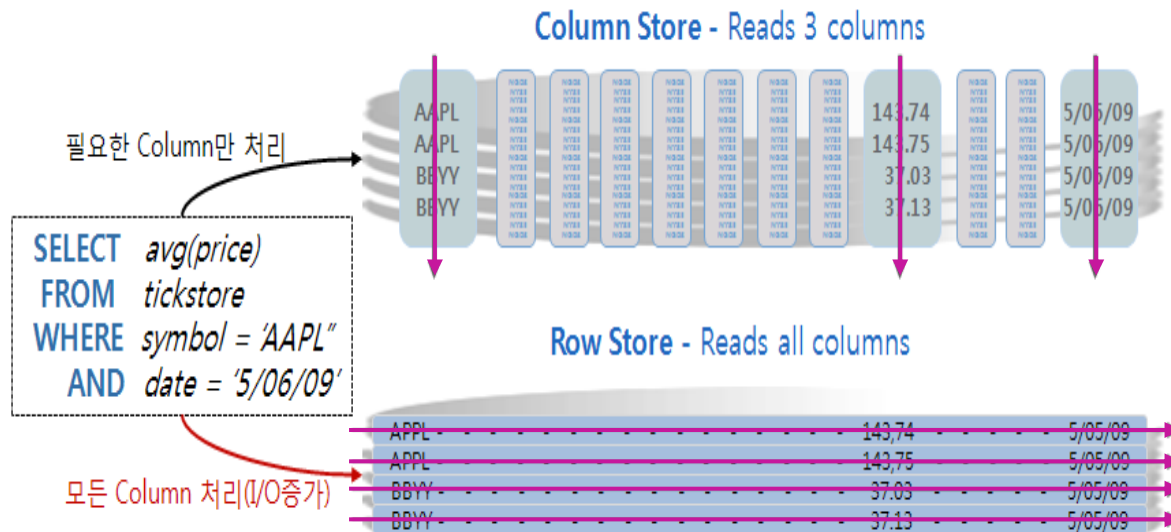
### VERTICA

- 동일 스펙/ 동일 구성의 서버
- 단순한 구성
- 아무 서버에나 작업을 요청하면 전 노드가 병렬 수행
- 시장의 신뢰도를 확보한 서버 사용
- 관리 및 사용 용이성 확보



# Vertica 아키텍처의 특징 > Native Columnar

대용량 데이터를 관리하는 DW 데이터베이스의 성능은 I/O를 얼마나 줄일 수 있느냐에 달려 있습니다. Vertica의 Columnar DBMS 아키텍처는 Query 수행에 필요한 Column 만을 읽어 올 수 있도록 설계되었기 때문에, Query 시 마다 모든 열을 읽어와야 하는 row 기반 DBMS와 비교하여 I/O 발생량을 획기적으로 감소시킬 수 있습니다.

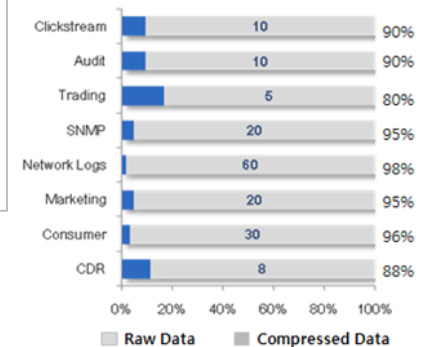
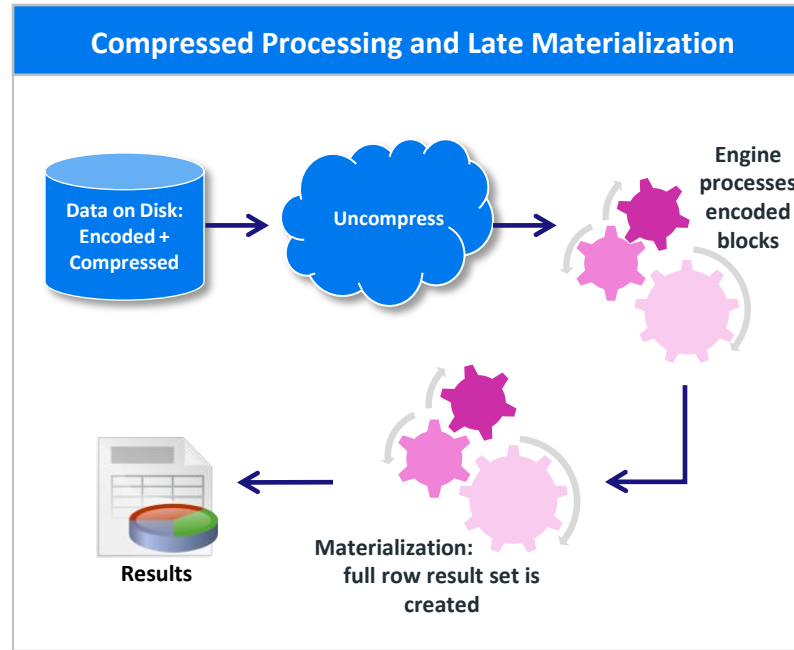
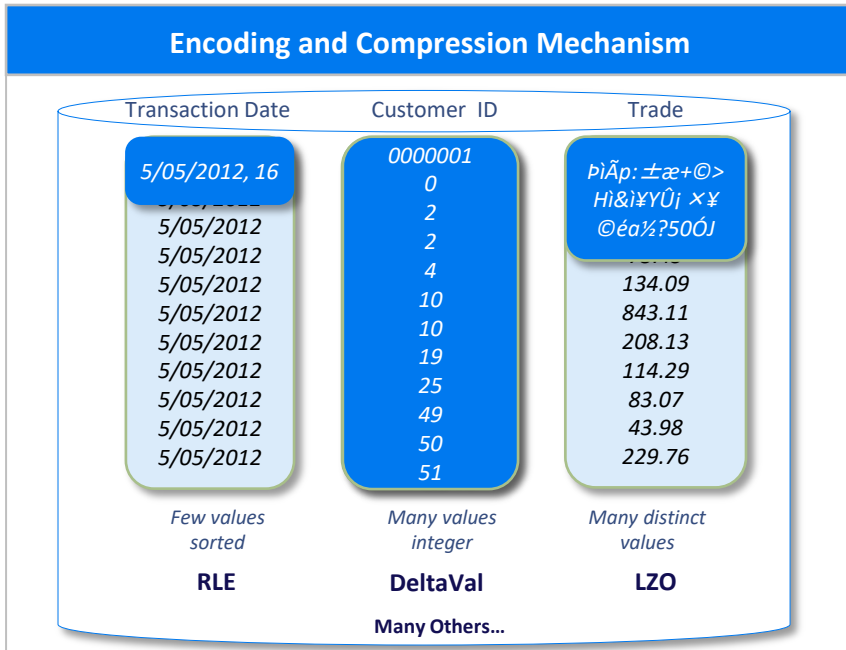


- 압축과 질의가 물리적인 I/O 레벨부터 컬럼 단위로 처리
- 컬럼 기반 저장 기술에 맞는 쿼리 옵티마이저
- 컬럼 단위 저장, 처리를 위한 별도의 옵션이나 절차가 불필요
- 컬럼 저장 구조에 최적화된 데이터 적재와 트랜잭션 처리
- 적은 하드웨어 리소스로 다른 DBMS와 동일한 작업 수행



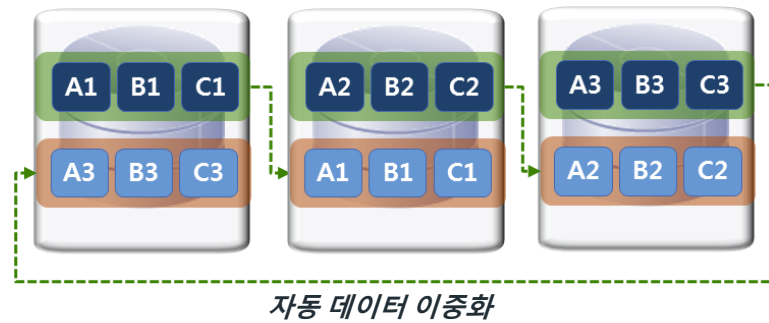
# Vertica 아키텍처의 특징 > 압축 및 인코딩으로 I/O 최소화

서로 다른 데이터 타입이 혼재되어 있어 압축률이 좋지 않은 row 기반 DBMS와는 달리, **동일한 데이터 타입을 가지는 column 단위로 데이터를 저장하는 column 기반 DBMS는 높은 압축율을 제공합니다.** Vertica에 내장된 12가지 데이터 인코딩 및 압축 알고리즘은 90% 이상의 압축율을 제공하여 스토리지 사용량을 효과적으로 절감할 수 있도록 합니다.



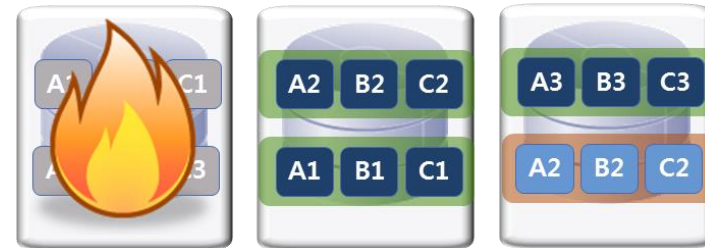
# Vertica 아키텍처의 특징 > 데이터 이중화로 무중단 서비스 제공

Vertica 는 저장 데이터의 이중화 기능을 이용하여 성능 향상과 함께 **노드 장애 시에도 중단 없는 서비스**를 가능하게 합니다. 데이터베이스 용량 확장을 위한 노드 추가 시나 유지보수를 위한 노드 제거 시에도 서비스 중단 없이 작업이 가능합니다.



데이터가 이중화되어 있어 서비스는 중단없이 지속  
시스템 장애 복구시 자동으로 클러스터 내의 다른 서버로 부터 데이터 동기화 수행

장애발생



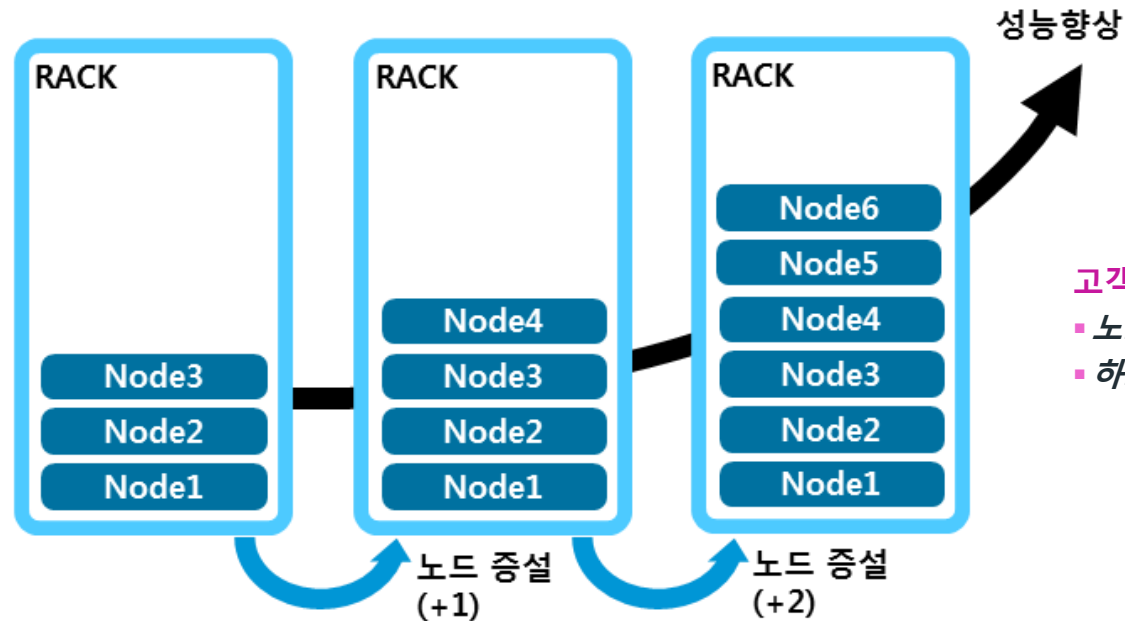
- RAID 기능과 유사한 데이터 이중화에 의한 노드 장애 무중단 지원
- 데이터 베이스 확장을 위한 노드 증설 시에도 서비스 무중단
- H/W 유지 보수를 위한 노드 제거 시에도 서비스 무중단 - (CPU/Memory/스토리지/OS 등)
- 스토리지 장애 시에도 해당 노드 무중단 (RAID 적용)
- 시스템 장애 복구시 자동으로 클러스터 내의 다른 서버로 부터 데이터 동기화 수행



노드 장애가 해결되어 다시 정상적으로 부트되면, 자동적으로 그동안 다른 노드에서 변경되었던 데이터가 중단되었던 노드로 동기화 됩니다. 이 과정 역시 온라인으로 자동 수행되므로 서비스의 중단은 발생하지 않습니다.

# Vertica 아키텍처의 특징 > 증설

데이터 및 사용자 증가 시 **노드 단위의 증설**을 통해 성능을 향상시키고, **온라인 데이터 재분배** 과정으로 데이터베이스 다운타임 없이 확장 가능합니다.



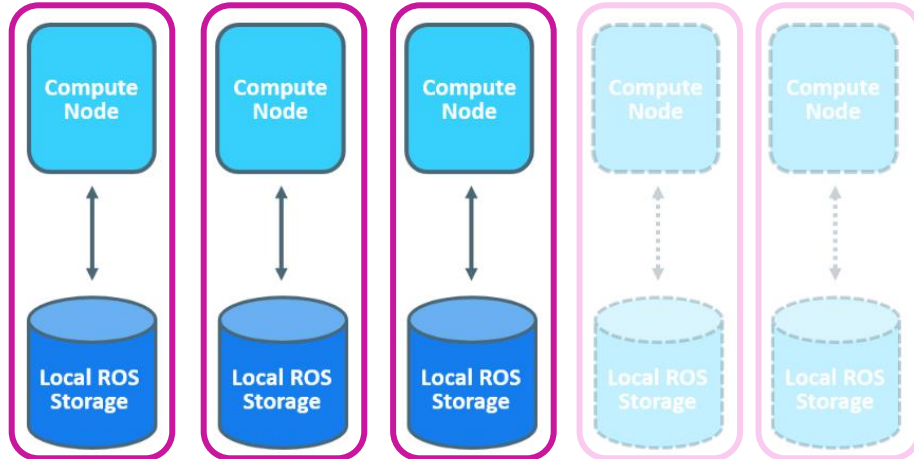
## 고객투자 보호

- 노드 단위 하드웨어 증설로 비용 최적화
- 하드웨어 사양에 관계없이 1TB 단위로 DB 라이선스 증설

# Vertica 아키텍처의 특징 > Eon Mode

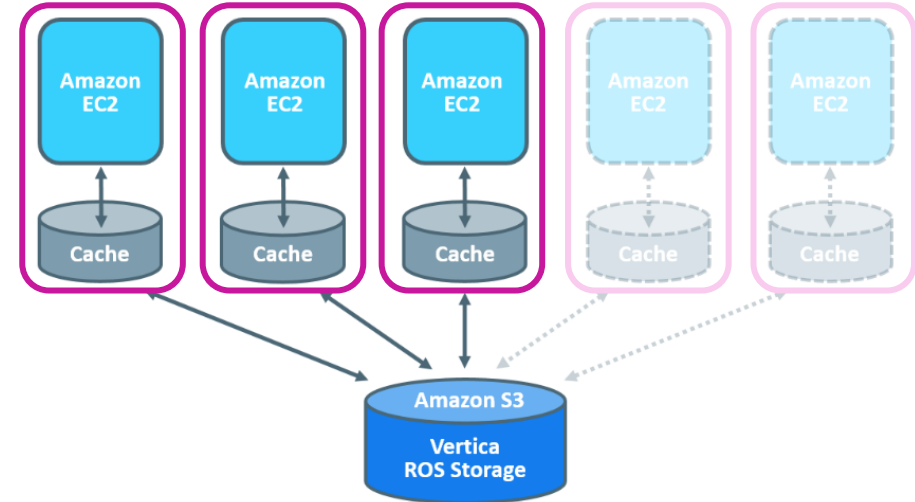
Vertica Enterprise Mode 외에 Object Storage(S3 compatible) 를 활용하여 **cs 분리를 지원하는 새로운 아키텍처가 지원됩니다.**

**Vertica Enterprise Mode**  
(On-premises, Cloud, 또는 Hybrid)



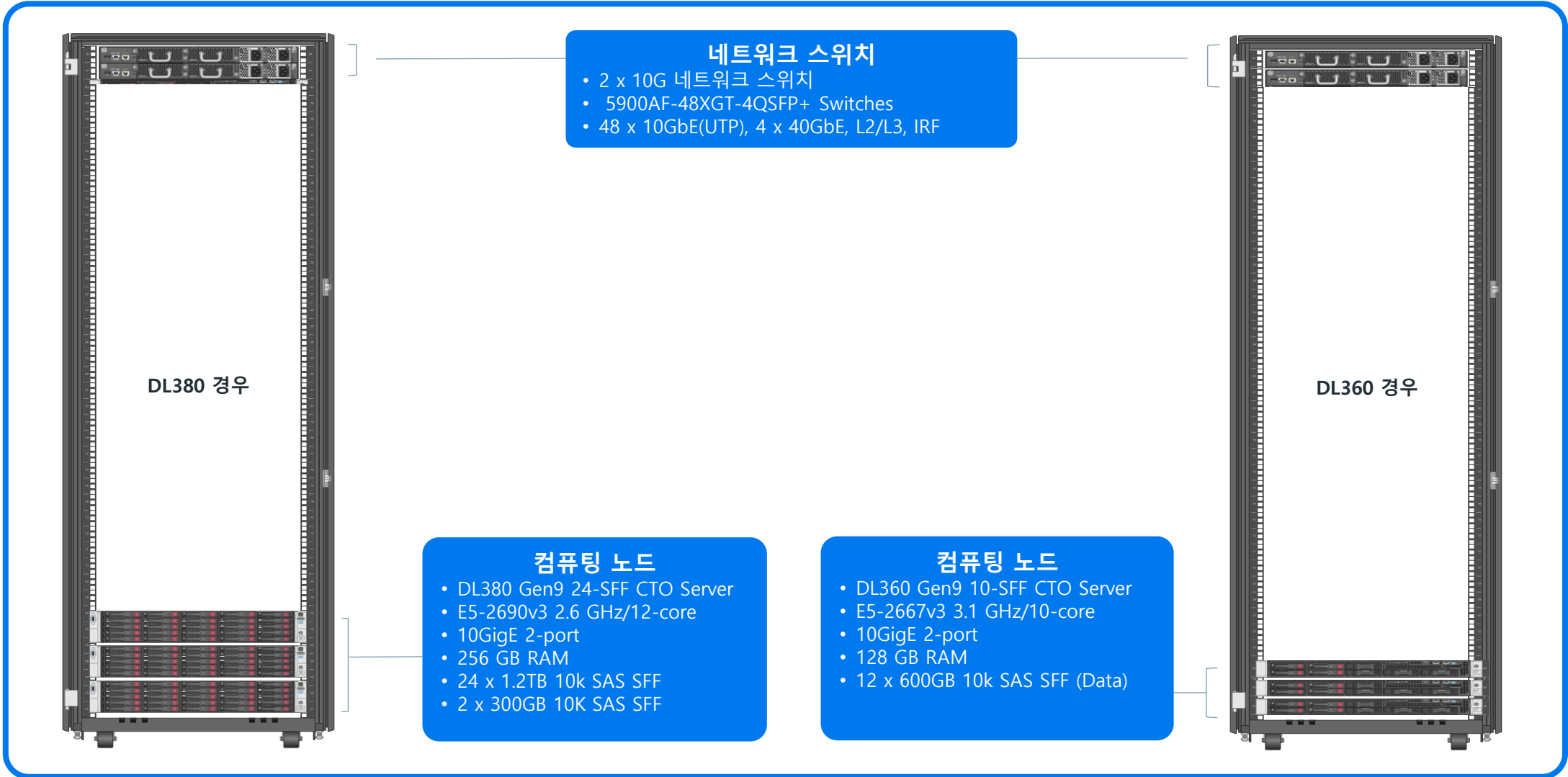
컴퓨팅 엔진과 스토리지가 강하게 결합되어 예측 가능한 워크로드를 원하는 기대 시간 내에 빠르게 처리하기 위한 아키텍처

**Vertica Eon Mode**  
(Amazon Web Services, On-Premise with Pure Storage)



클라우드 이코노믹스의 동적 워크로드 요구 사항에 대응하는 컴퓨팅 리소스만 독립적으로 확장이 가능한 아키텍처

# Standard Vertica System(예시)





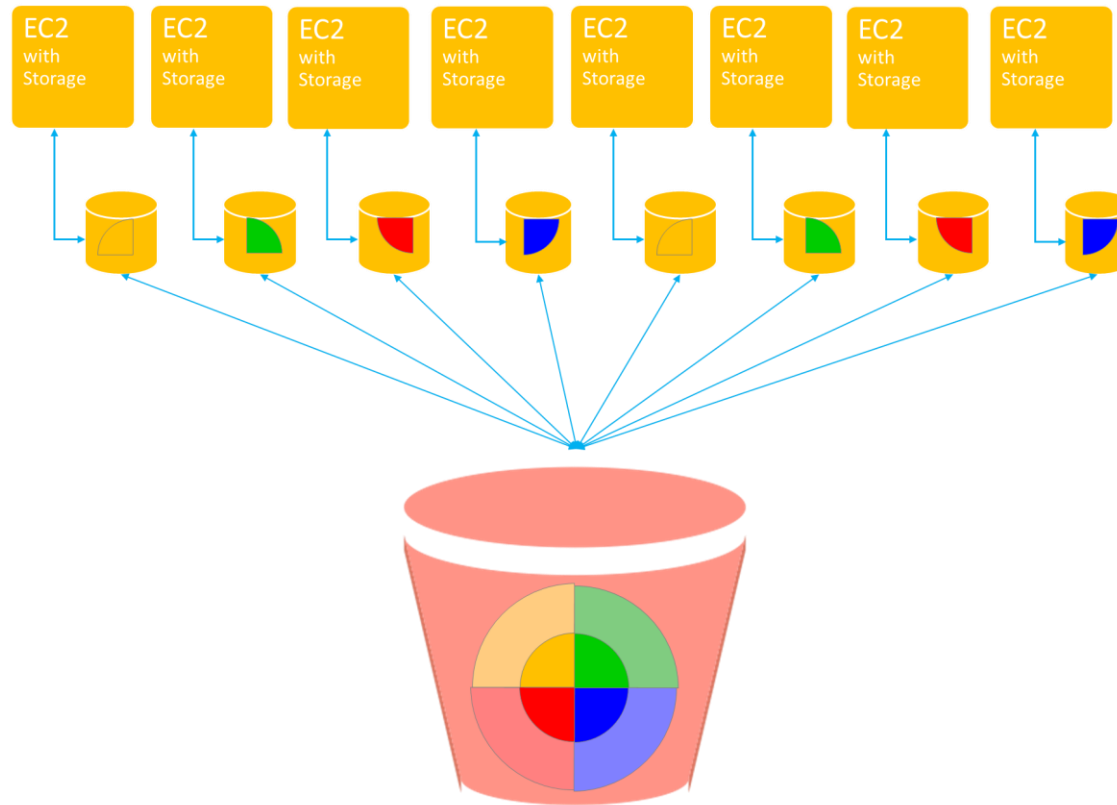


VERTICA

# Vertica Eon Mode

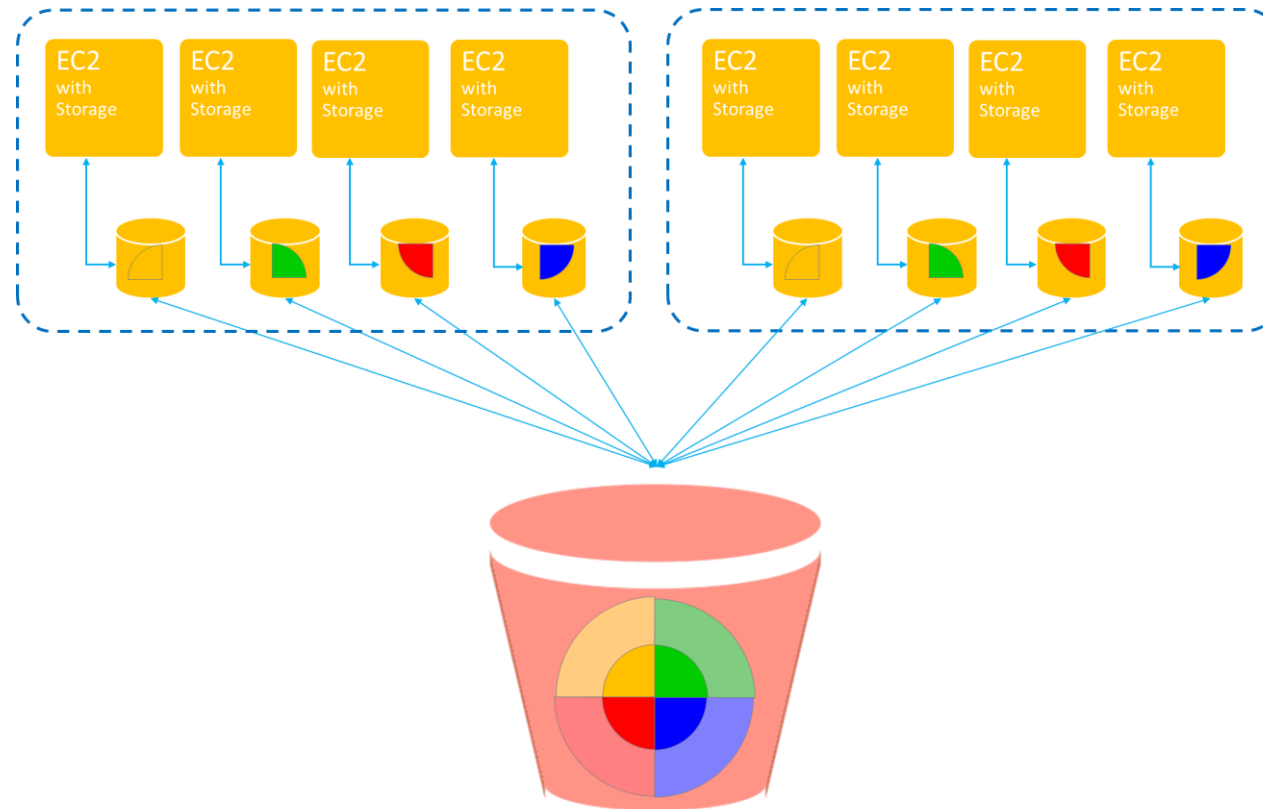
# 쉽고 빠른 확장성

동시성이나 사용자가 많아지면 바로 Scale Out 을 수행할 수 있으며, 영구 데이터가 별도의 공간에 있어 **Data Rebalancing 이 불필요**  
인스턴스는 기존에 만든 VM 이미지로 빠르게 생성하여 추가할 수 있으며, 일시적인 증가라면 해당 작업을 마친 후 Scale In 을 수행



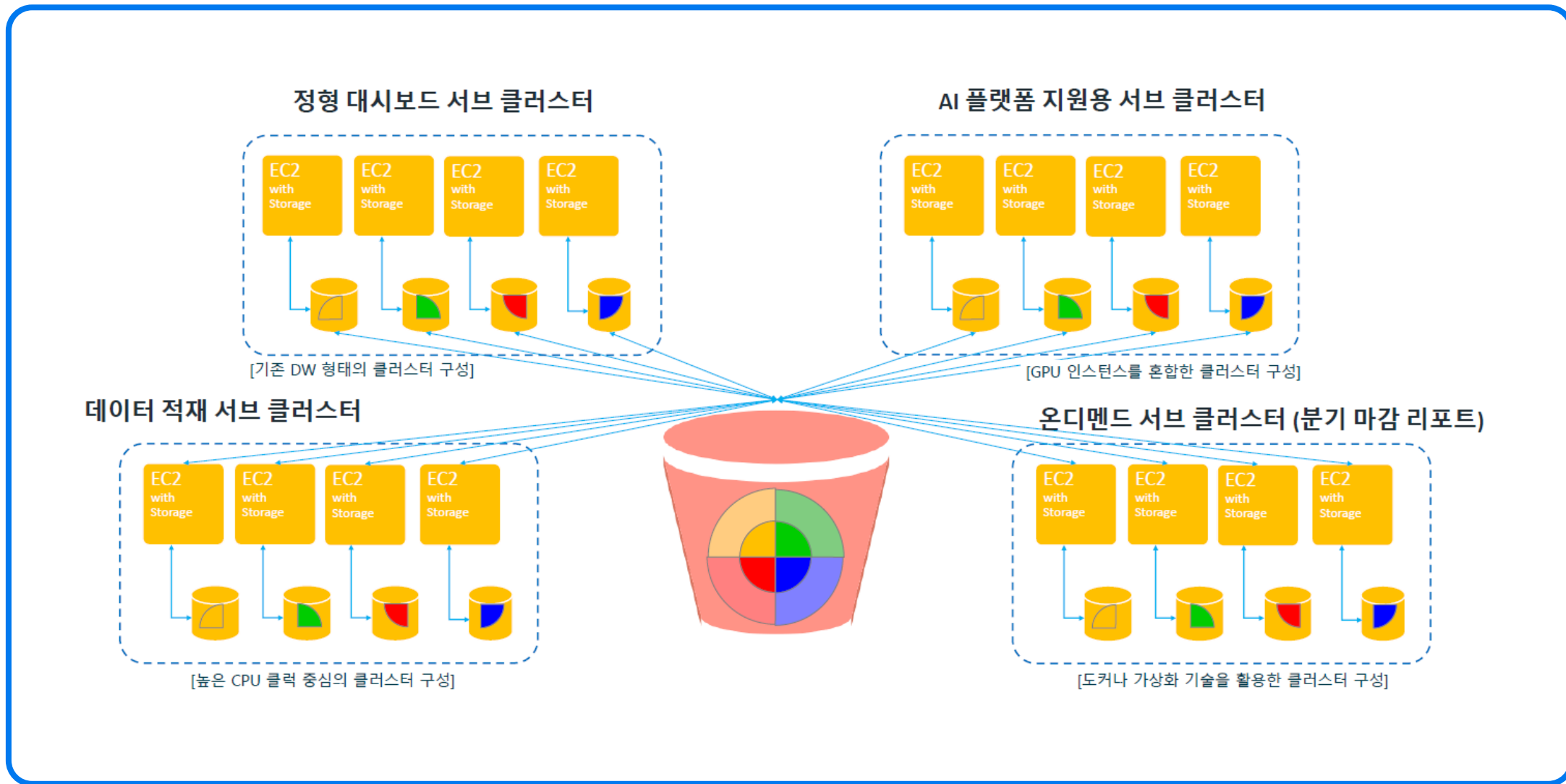
# 서브클러스터를 활용한 워크로드 분리

추가적으로 요구되는 업무의 특징이 기존의 업무와 다른 성격을 갖고 있는 경우에는, 서브 클러스터를 구성하는 것이 효과적  
클러스터별로 다른 형태의 서버 구성 클럭, 메모리, GPU 등 클러스터별로 다른 형태의 업무 수행 적재, 배치, 정형업무, 머신러닝 등





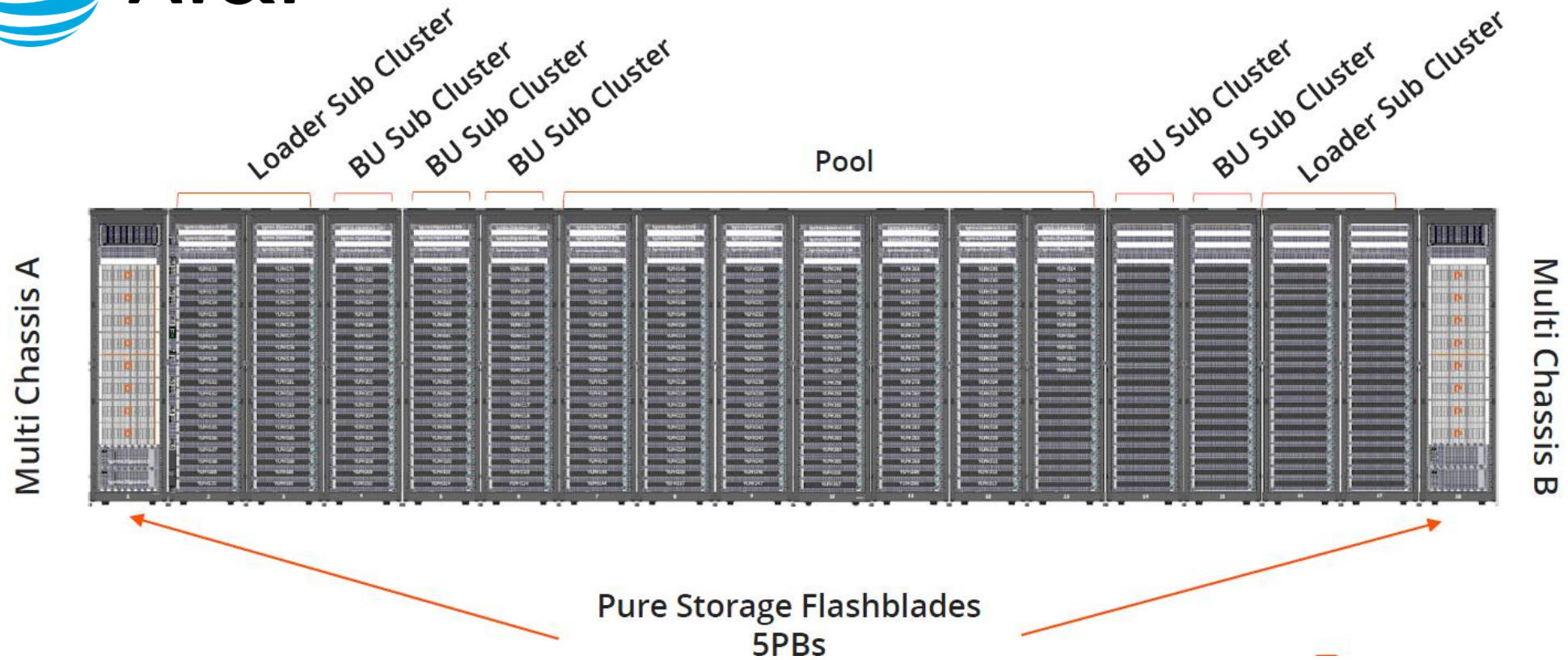
# 서브 클러스터를 활용한 워크로드 분리 활용 예



# Physical deployment



AT&T





VERTICA

# In-DB Machine Learning



# 다양한 머신 러닝 알고리즘 지원

순수 내재화 된 SQL 기반 함수로 지원



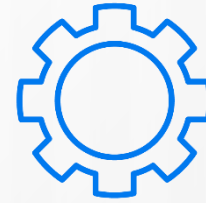
SQL 데이터베이스

+



고급분석과 머신러닝

+



쿼리 엔진

Data Analysis

Data Preparation

Modeling

Evaluation

Deployment



Linear  
Regression



Logistic  
Regression



K-Means  
Clustering



Naive  
Bayes



Support Vector  
Machines



Random  
Forrest

# End-to-end 전체 머신러닝 주기를 모두 지원

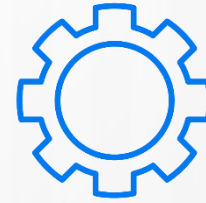
버티카만으로 머신러닝 업무 수행이 가능



SQL 데이터베이스



고급분석과 머신러닝



쿼리 엔진

## Data Analysis

- Statistical Summary
- Time Series
- Sessionize
- Pattern Matching
- Date/Time Algebra
- Window Partition
- Sequences
- And more...

## Data Preparation

- Outlier Detection
- Normalization
- Imbalanced Data Processing
- Sampling
- Missing Value Imputation
- And More...

## Modeling

- SVM
- Random Forests
- Logistic Regression
- Linear Regression
- Ridge Regression
- Naïve Bayes
- Cross Validation
- And More...

## Evaluation

- Model-level Stats
- ROC Tables
- Error Rate
- Lift Table
- Confusion Matrix
- R-Squared
- MSE
- And More...

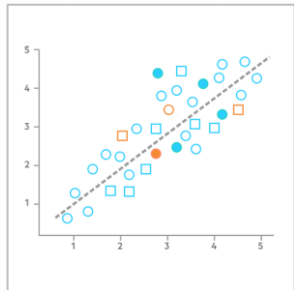
## Deployment

- Deploy Anywhere
- In Database Scoring
- Massively Parallel Processing
- Speed
- Scale
- Security
- And More...

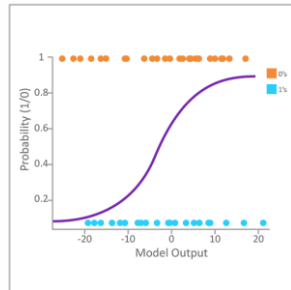
# 단순한 SQL 호출로 수행

데이터베이스에서 인식하고 있는 데이터에 대해 SQL로 분석 함수 호출

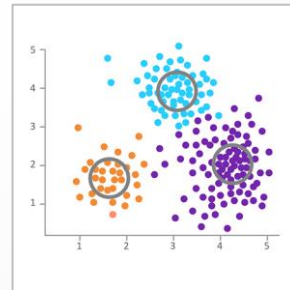
SQL



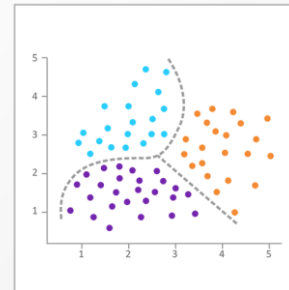
Linear  
Regression



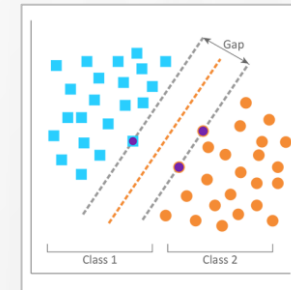
Logistic  
Regression



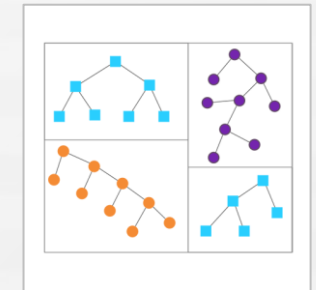
K-Means  
Clustering



Naive  
Bayes



Support Vector  
Machines



Random  
Forrest

버티카에서 사용자는 모델을 생성하고 학습하고 배포하는 것이 가능

# 머신 러닝 모델 생성

단순한 SQL 함수로 수행

Creates new model

Select table/view that contains training data

Select column with dependent variable

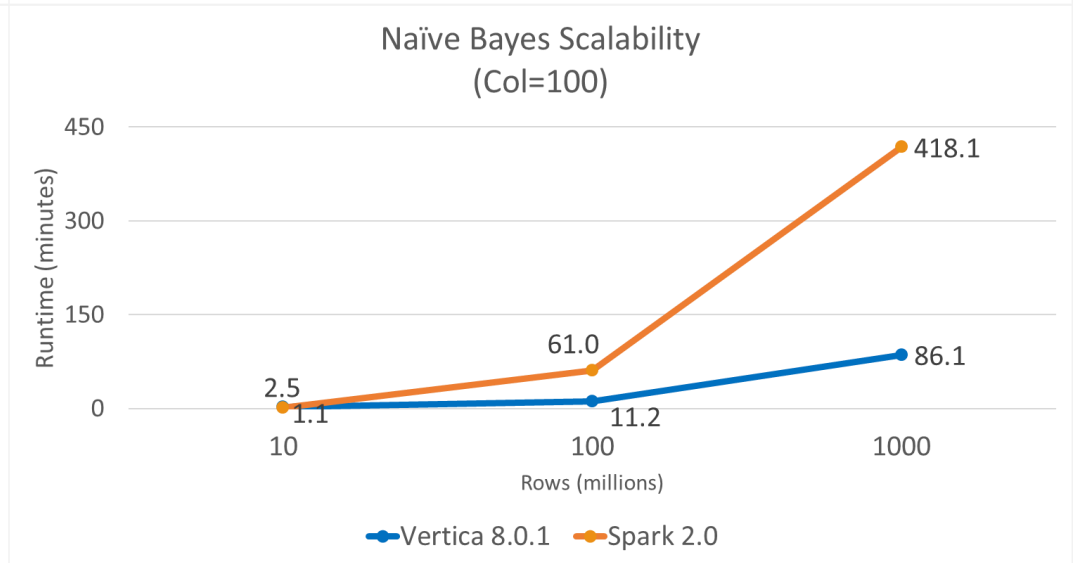
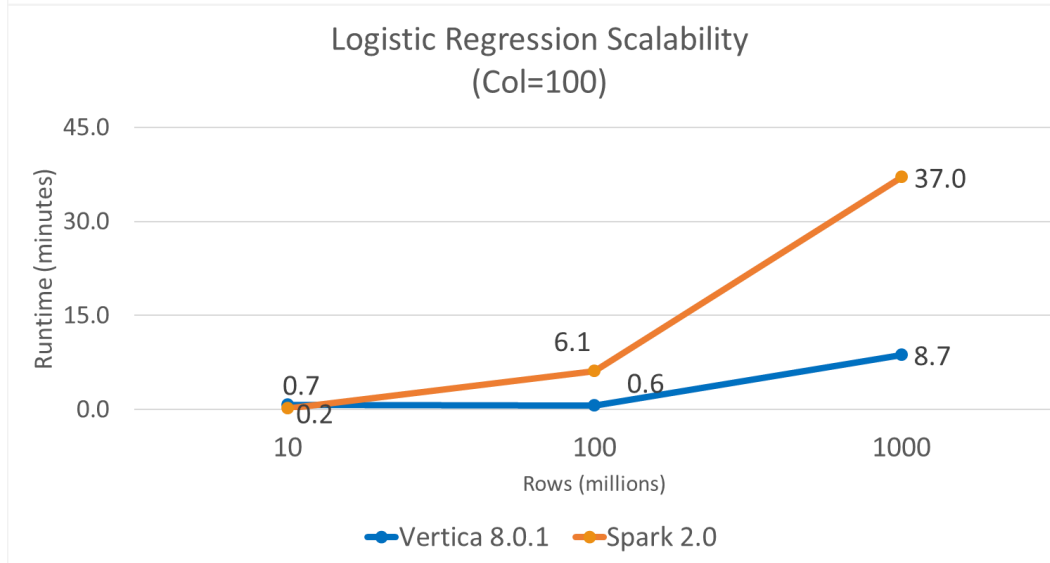
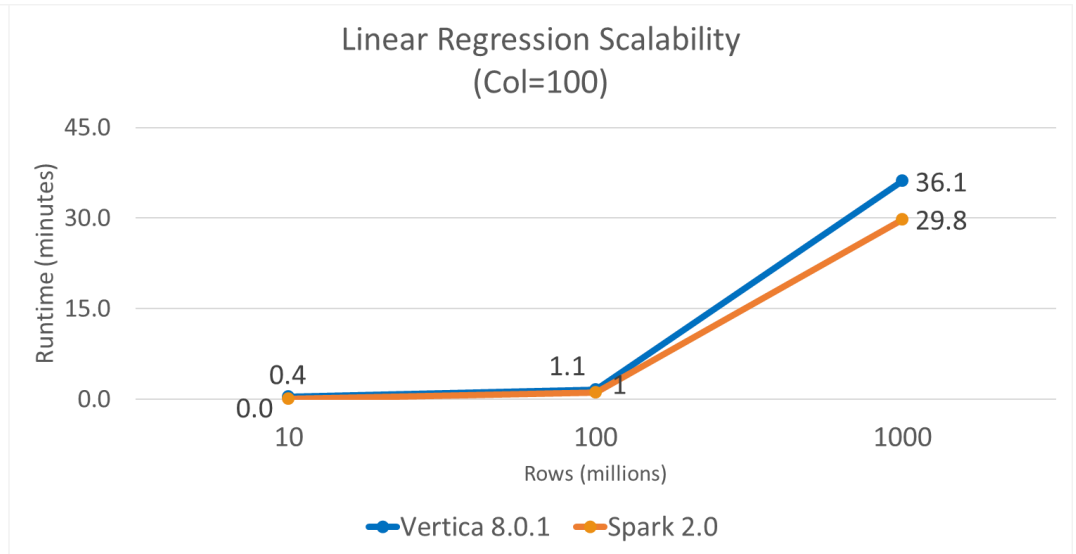
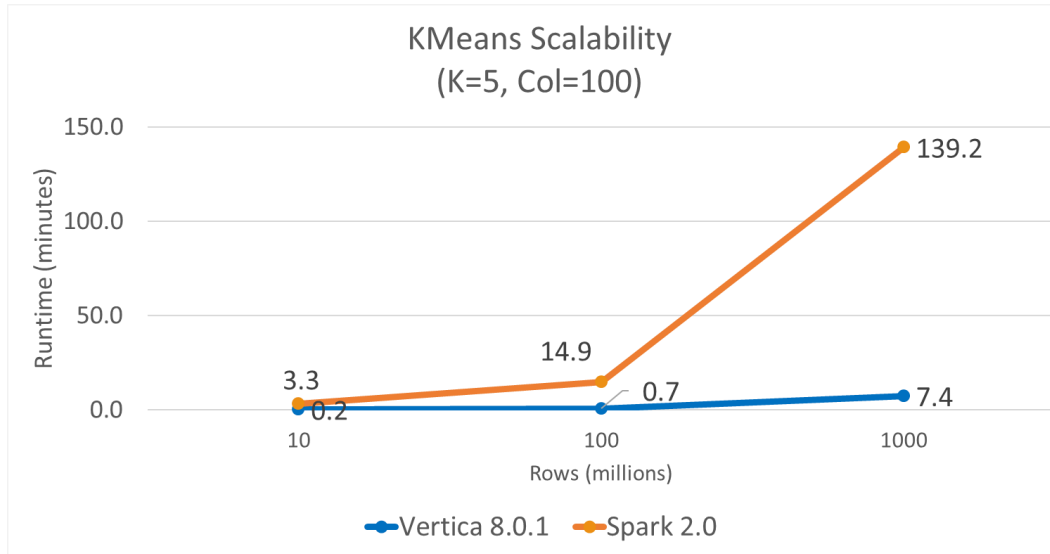
Select columns with independent variables

```
LINEAR_REG ( 'model_name', 'input_relation', 'response_column', 'predictor_columns'
  [ USING PARAMETERS [exclude_columns='col1, col2, ... coln',]
    [optimizer='value',]
    [epsilon=value,]
    [max_iterations=value,]
    [regularization= 'value',]
    [lambda= value,]
    [alpha = value]])
```

Optional parameters  
for model building



# Spark 와의 동시성 성능 비교

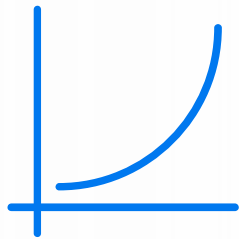


# 버티카 In-DB 머신러닝 특징점

다양한 방법으로 데이터 분석 업무를 지원

## 확장성

병렬 처리 가능한  
데이터 분석



쉬운 SQL 문법으로  
더 많은 사용자들이  
더 많은 데이터에 대해  
머신러닝을 수행 가능

## 고성능

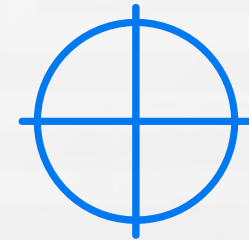
머신 러닝 요건에 대한  
빠른 비즈니스 대응



버티카의 병렬 처리 기능을  
머신러닝에도 적용하여  
빠른 성능을 보장

## 정확성

지속 가능하게 정확도를  
높이는 반복 학습 수행



샘플링 된 데이터가 아닌  
전수 데이터에 대해  
지속적인 학습 수행으로  
정확도 높은 빅데이터 분석 달성

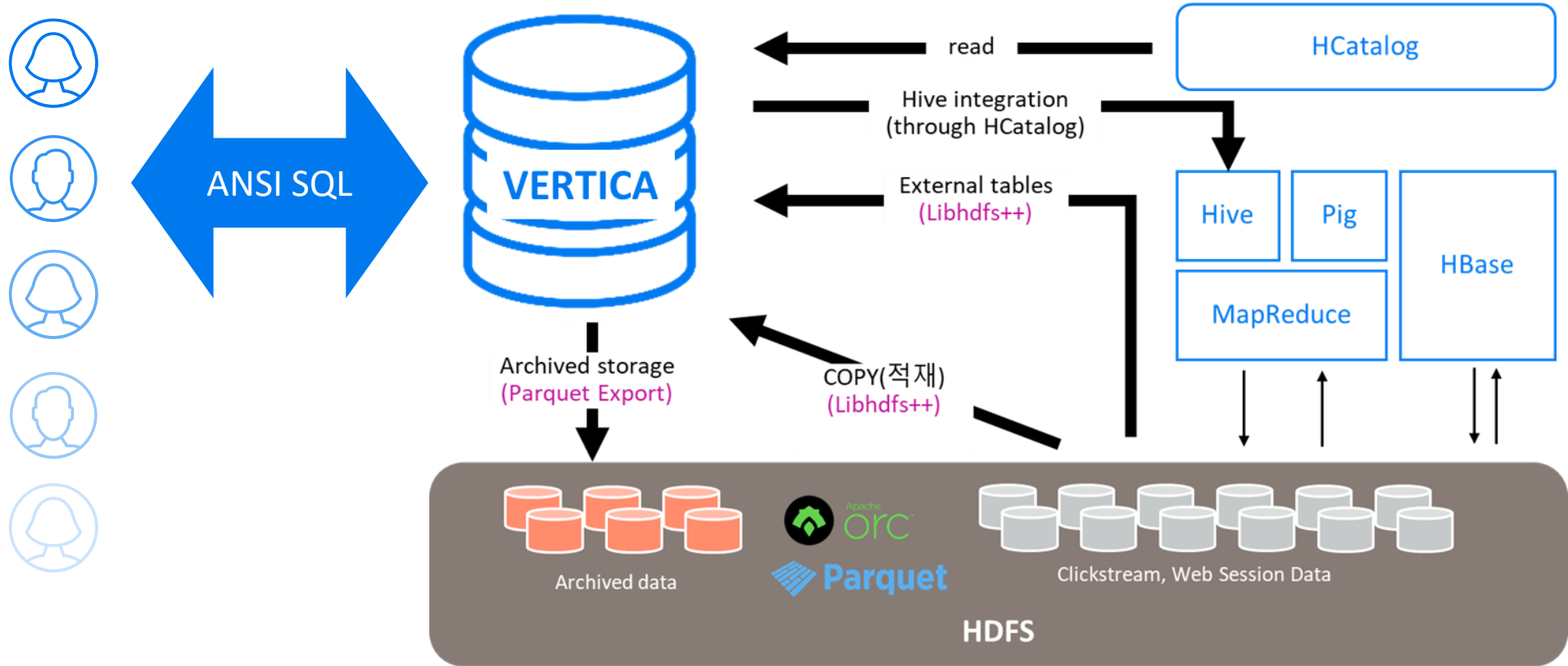


VERTICA

# Integrating with Hadoop

# Integration Points

추가적인 장비나 별도의 솔루션 설치 없이 하둡 연계를 지원하여 버티카를 통해 **dw와 하둡 데이터 연계 분석**을 지원

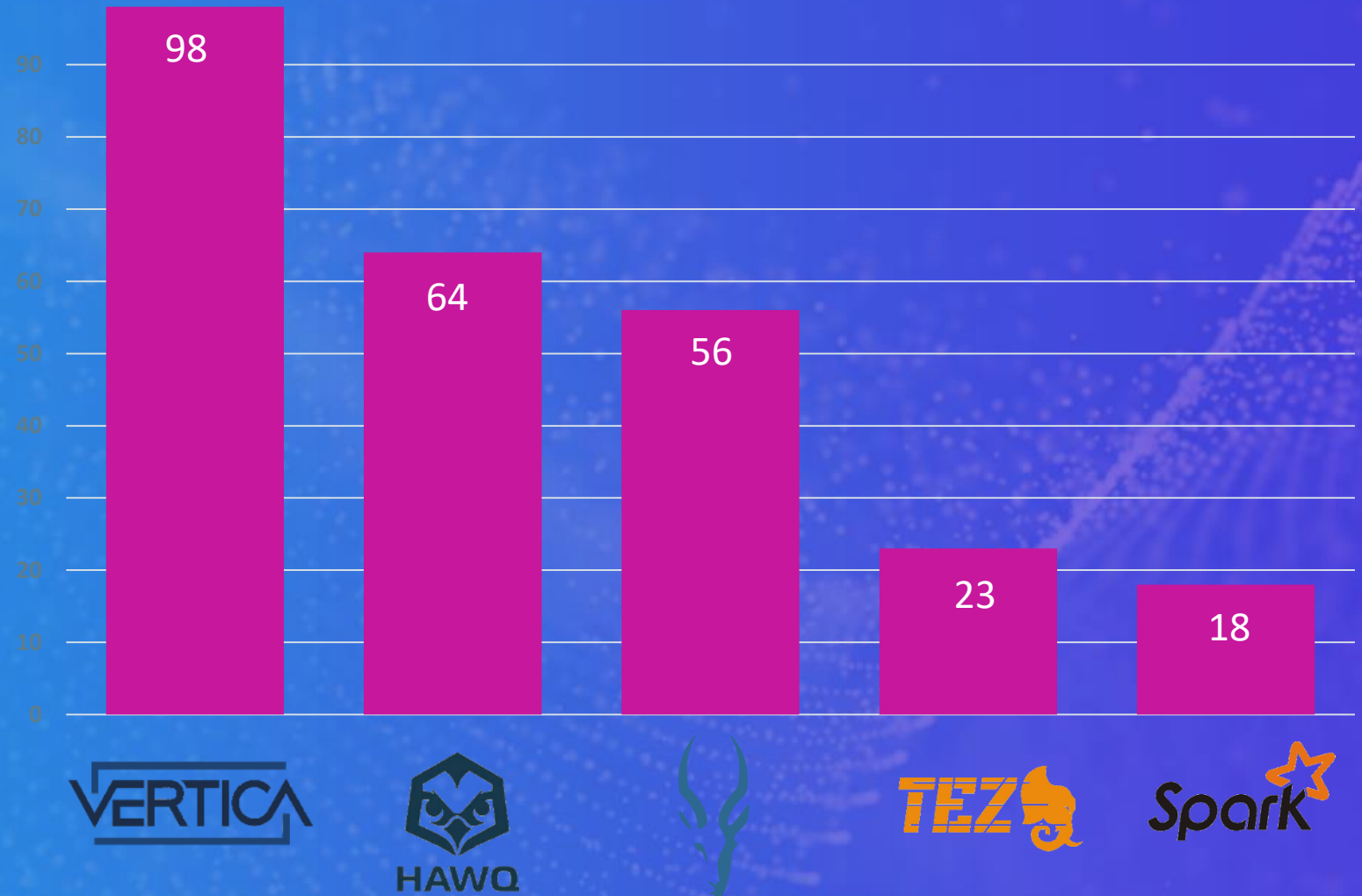




# SQL on Hadoop 기술 비교

- 분석 시스템 비교를 위한 표준 Benchmark 방법론인 TPC-DS 기준
- TPC-DS 99개의 쿼리중 각 솔루션별로 성공한 개수
- 표준 SQL에 대한 호환성에 문제가 있는 솔루션의 경우 추가 개발 공수 필요
- HAWQ 기반의 SQL은 대부분 쉽게 이식될 수 있음

Successful Unaltered TPC-DS Queries



# Vertica Enterprise vs. VSOH(Hortonworks) vs. TEZ(Hortonworks)

- Vertica Enterprise는 Hive on Tez 대비 약 14배의 빠른 성능을 보임
- ORC를 사용한 Vertica SQL on Hadoop은 Hive on Tez 대비 약 8배 빠른 성능을 보임
- Hive on Tez 는 전체 99개 쿼리 중 40개를 실패



Seconds to complete benchmarks  
(of runnable queries)



# VSOH 사례 - AT&T

500 x Hadoop Nodes (Hortonworks, ORC)



70 x VSOH Nodes



Kerberos Enabled



Tableau



DBViz

# VERTICA