**Microsoft**

# Data Warehousing and Big Data

Technology Deck

# Microsoft is a leader for...

CHALLENGERS · LEADERS · Oracle · Microsoft · Amazon Web Services · IBM · SAP · MongoDB · MariaDB · DataStax · EnterpriseDB · MarkLogic · InterSystems · Percona · Redis Labs · FairCom · MapR · Cloudera · Couchbase · Neo Technology · Fujitsu · Altibase · MemSQL · VoltDB · NuoDB · TmaxSoft · Clustrix · Actian · Basho Technologies · Aerospike · Hortonworks · Orient Technologies · McObject · NICHE PLAYERS · VISIONARIES · ABILITY TO EXECUTE · COMPLETENESS OF VISION · As of October 2015

CHALLENGERS · LEADERS · Tableau · Qlik · Microsoft · Birst · SAS · Alteryx · SAP · MicroStrategy · Domo · GoodData · Salesforce · Logi Analytics · IBM · ClearStory Data · Board International · Pentaho · Sisense · TIBCO Software · Pyramid Analytics · Information Builders · Yellowfin · BeyondCore · Platfora · Datawatch · NICHE PLAYERS · VISIONARIES · ABILITY TO EXECUTE · COMPLETENESS OF VISION · As of February 2016

CHALLENGERS · LEADERS · Oracle · Teradata · Microsoft · IBM · SAP · Amazon Web Services · HPE · 1010data · Infobright · MarkLogic · Exasol · Cloudera · Actian · MapR Technologies · Hortonworks · Transwarp · MongoDB · Pivotal · Kognitio · MemSQL · Hitachi · NICHE PLAYERS · VISIONARIES · ABILITY TO EXECUTE · COMPLETENESS OF VISION · As of February 2016

# Contents

SQL Server Data Warehouse Family

    SQL Server 2016

    APS Appliance
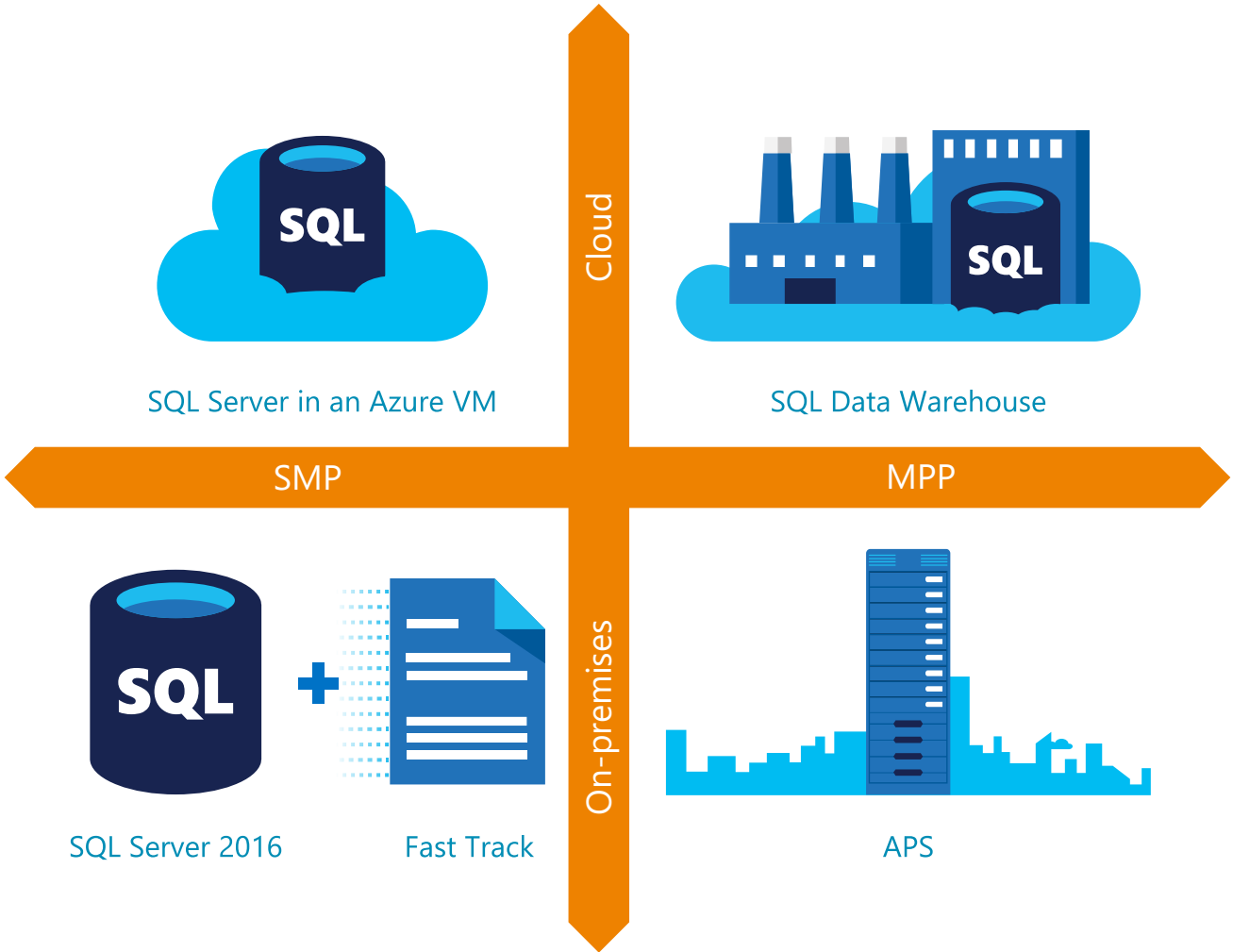
    SQL Data Warehouse

Microsoft Big Data Solutions

    Azure Data Lake

    HD Insight

    Azure Blob Storage

# SQL Server Data Warehousing solutions



SQL Server in an Azure VM

SQL Data Warehouse

SMP

MPP

Cloud

On-premises

SQL Server 2016      Fast Track

APS

## Symmetric multi-processing (SMP)

On-premises: SQL Server 2016 or SQL Server Fast Track Data Warehouse

Cloud: SQL Server in an Azure VM

## Massively parallel processing (MPP)

On-premise: Analytics Platform System (APS)

Cloud: Azure SQL Data Warehouse

# SQL Server 2016
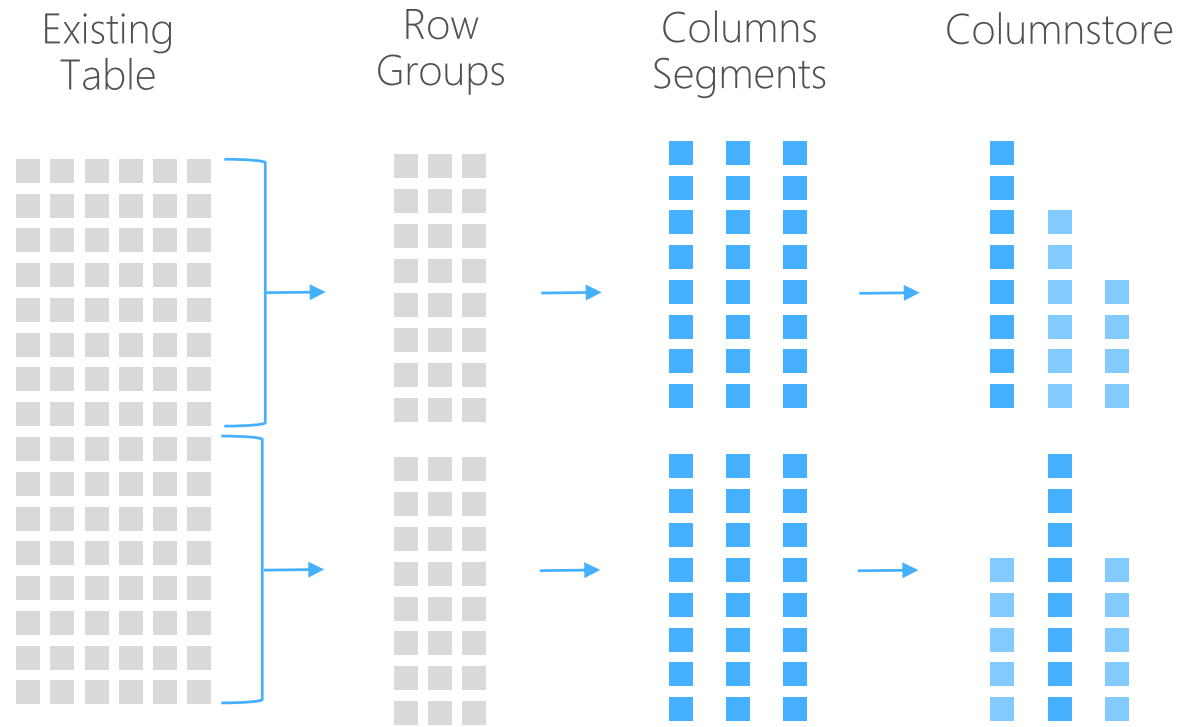
# Columnstore: Query performance and data compression

Existing
Table

Row
Groups
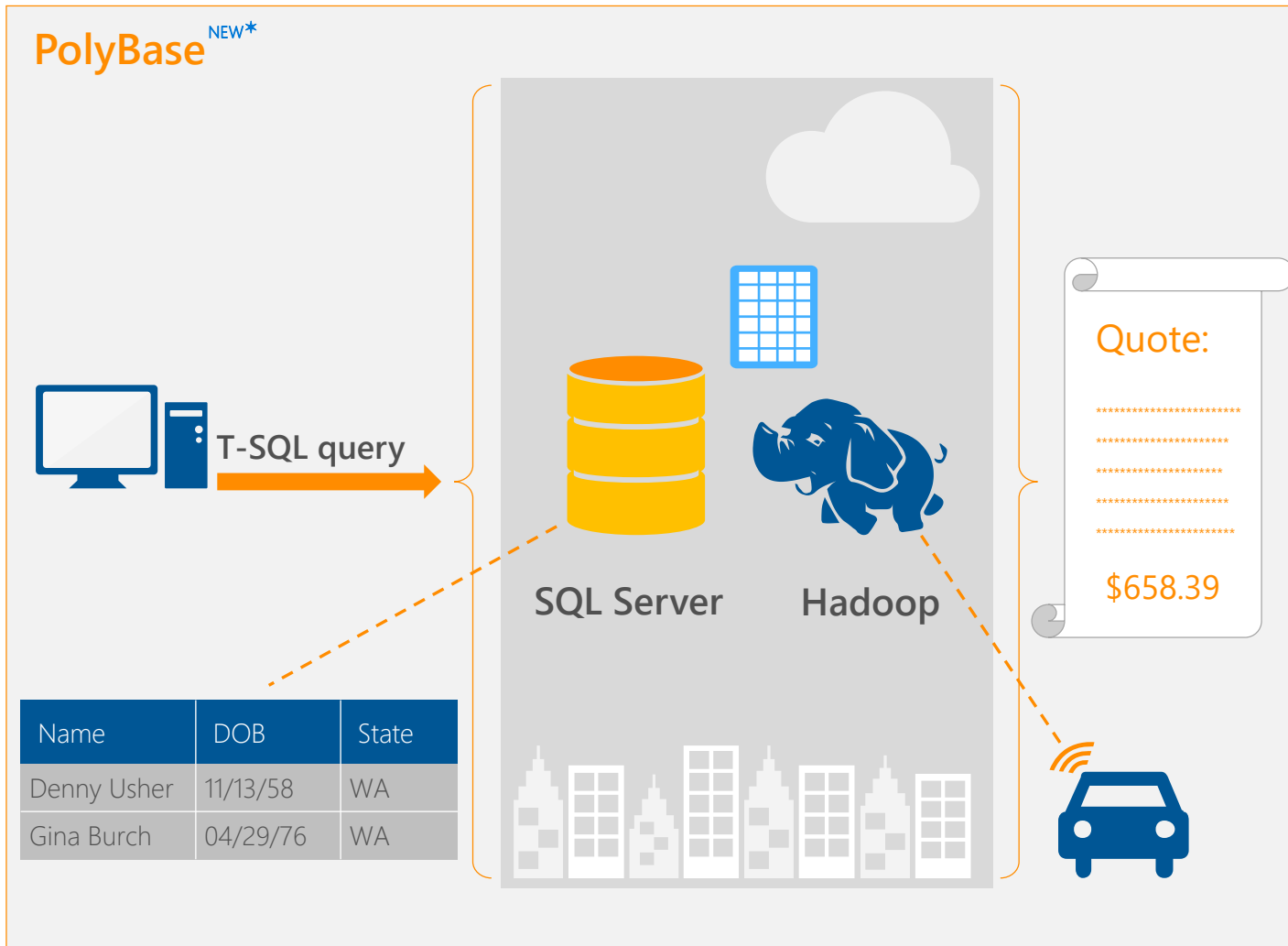
Columns
Segments

Columnstore

Columnstores are data structures organized in a column-based manner (as opposed to a row-based, traditional table)

Effective in scenarios where indexed columns have several repeated values

Appropriately designed columnstore indexes yield up to 100x the query performance and 10x the data compression of a traditional rowstore (table)

Compressed column segments are added to the columnstore.

# Remove the complexity of big data
## T-SQL over Hadoop

**PolyBase** NEW*

T-SQL query

| Name | DOB | State |
|------|------|-------|
| Denny Usher | 11/13/58 | WA |
| Gina Burch | 04/29/76 | WA |

SQL Server    Hadoop

Quote:

***************************
***************************
***************************
***************************
***************************

$658.39

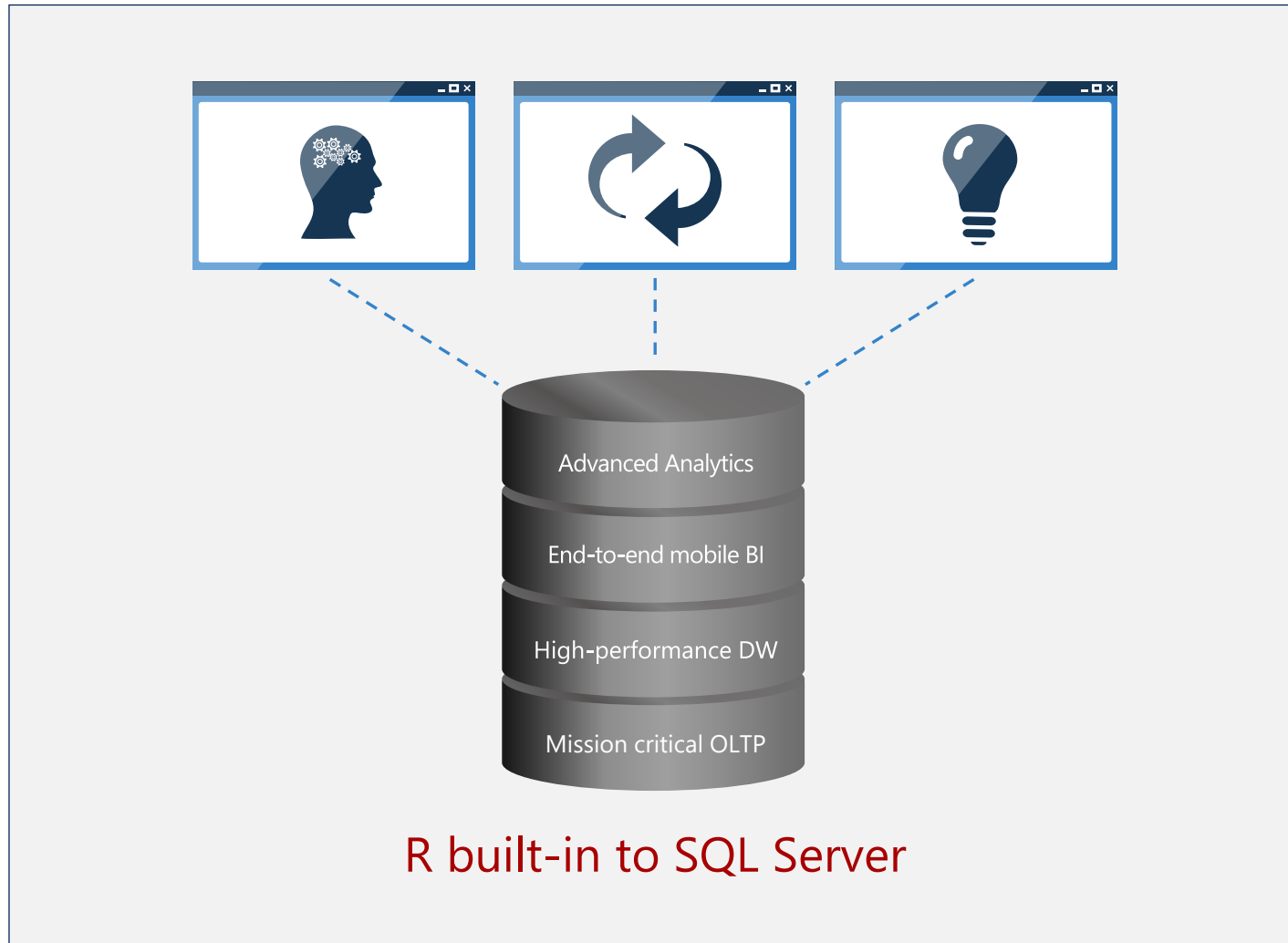Manage structured & unstructured data

**Simple T-SQL** NEW*
to query Hadoop data (HDFS)

**JSON support** NEW*

# In-database Advanced Analytics
Build intelligent applications with SQL Server R Services

Advanced Analytics

End-to-end mobile BI

High-performance DW

Mission critical OLTP

**R built-in to SQL Server**

**R built-in to your T-SQL** NEW*

**Real-time operational analytics** NEW*
without moving the data

**Open Source R with in-memory & massive scale** NEW* - multi-threading and massive parallel processing

# SQL Server 2016 (SMP) Reference Architectures

## Azure Virtual Machine Image for SQL Server Data Warehouse

SQL Server 2016 pre-built VM image in the Azure gallery

Disk Configuration for Data Warehousing
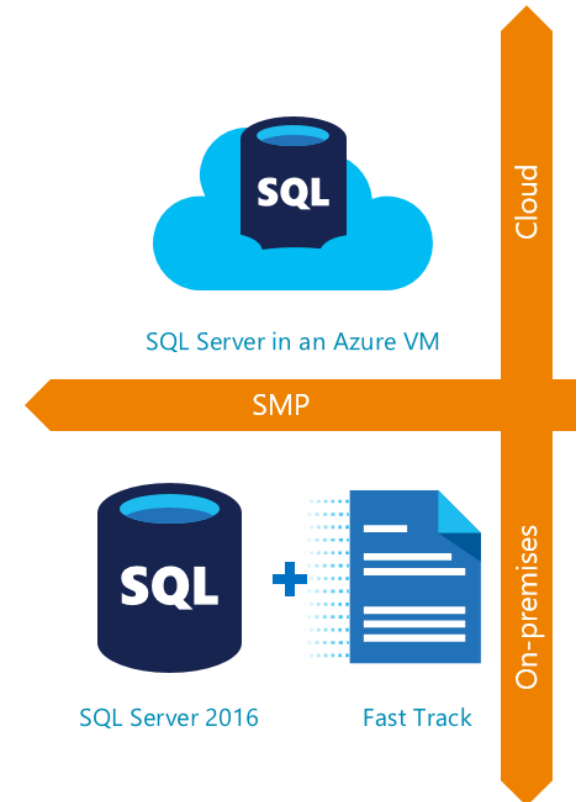
Developer Edition, BYOL, or per-hour Billing

Bottomless storage with Azure Blob Storage of Database files or Polybase

## Data Warehouse Fast Track

On-Prem Reference Architecture Implementations

HP, Dell, Lenovo, and other vendors

Tested Configurations from 5TB to 200TB

SQL Server in an Azure VM

Cloud

SMP

On-premises

SQL Server 2016    Fast Track

# SQL Server 2016 MPP Solutions

## SQL Data Warehouse

Data Warehouse-as-a-service
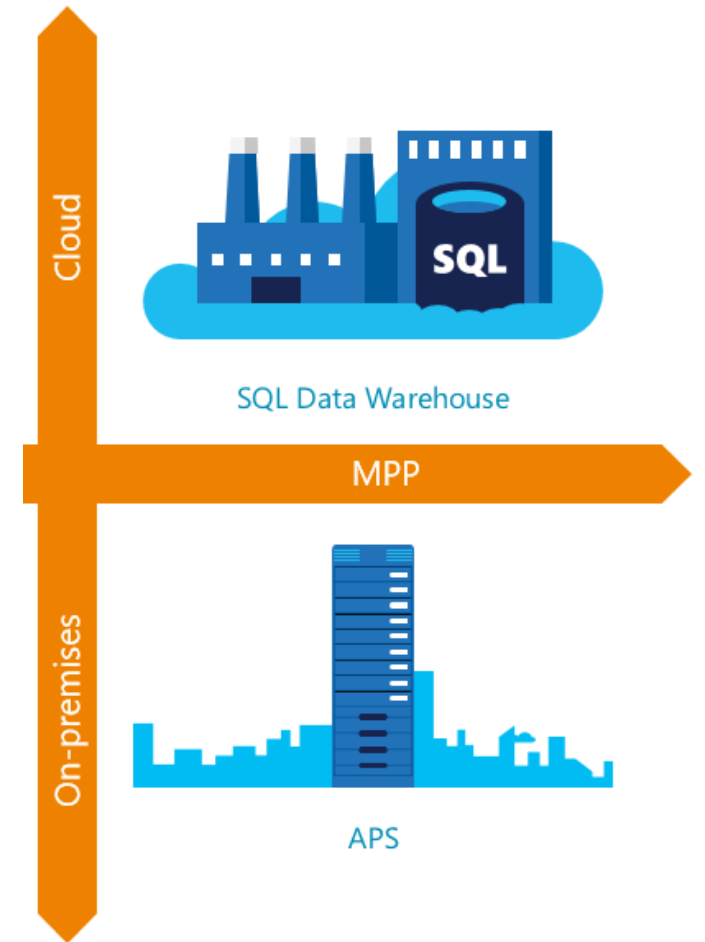
Elastic Scale in the Cloud

Polybase Connectivity to Azure Blob Storage

## Microsoft APS

On-Prem Data Warehouse Appliance

Partial-rack to multi-rack configurations

Polybase Connectivity to Azure Blob Storage and Hadoop

# Scaling out your data to petabytes

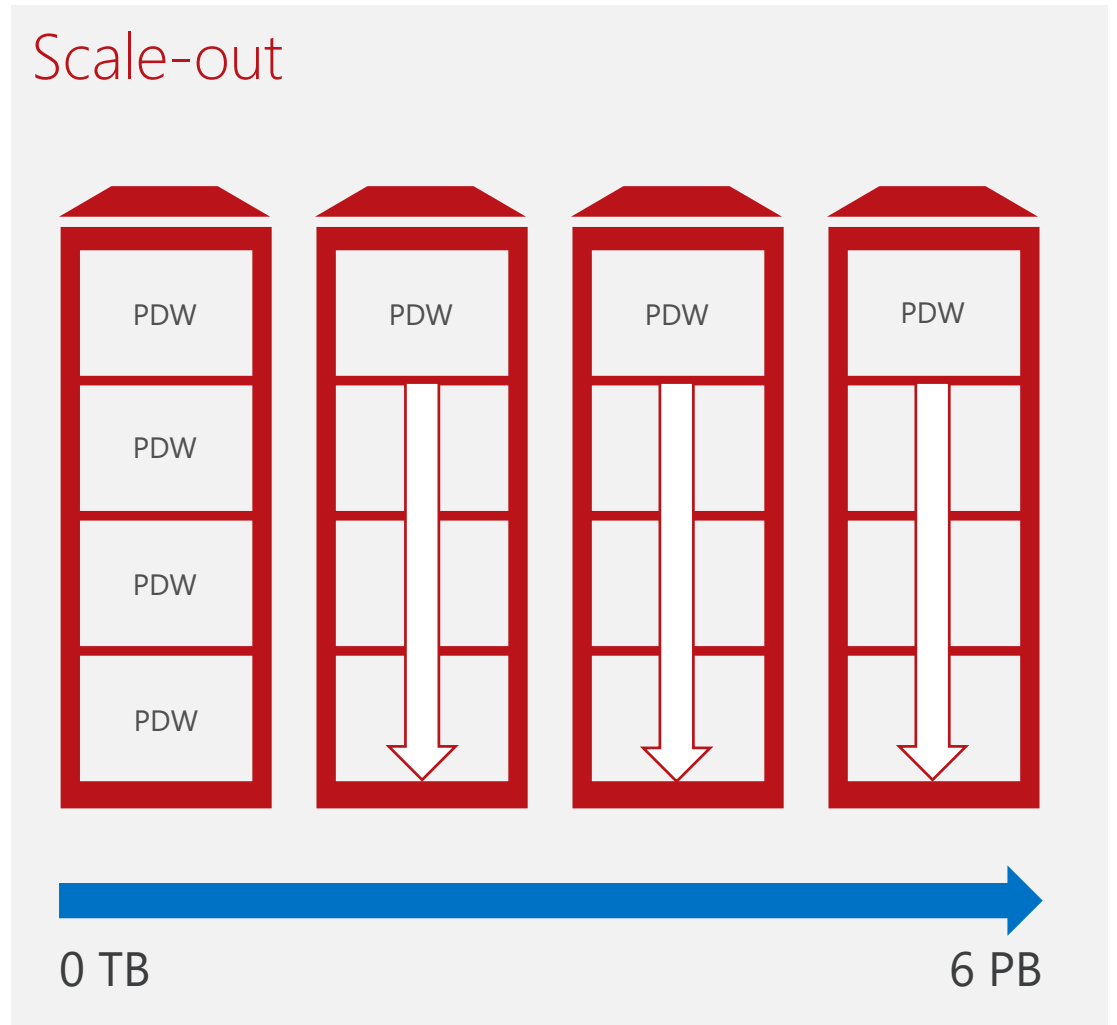## Scale-out technologies in Analytics Platform System

Multiple nodes with dedicated CPU, memory, and storage

Ability to incrementally add hardware for near-linear scale to multiple petabytes

Ability to handle query complexity and concurrency at scale

No "forklift" of prior warehouse to increase capacity

Ability to scale out PDW or Azure Blob Storage

Scale-out

PDW    PDW    PDW    PDW

PDW

PDW

PDW

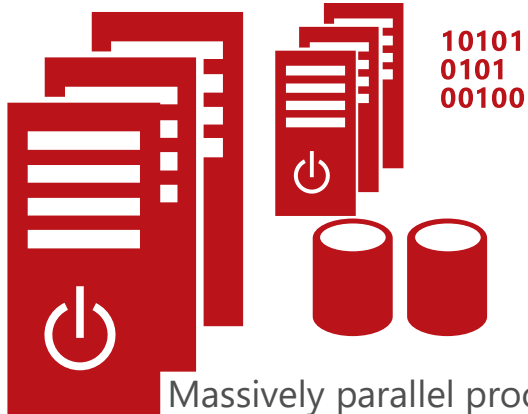0 TB                                          6 PB

# Azure SQL Data Warehouse

A relational data warehouse as a service, fully managed by Microsoft

Industry's first elastic cloud data warehouse with enterprise-grade capabilities

Support for your smallest to largest data storage needs while handling queries up to 100x faster

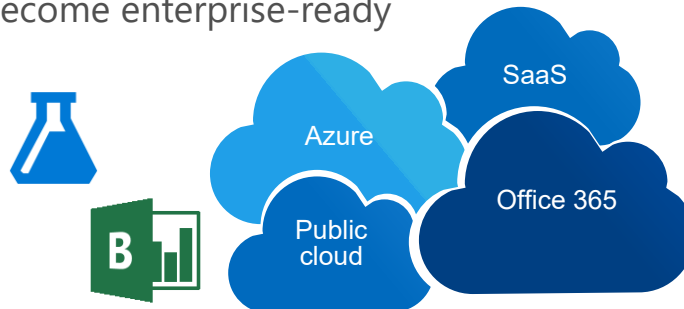## Elastic scale & performance

10101
0101
00100

Massively parallel processing

Scale to petabytes of data

Instant-on compute scales in seconds

Query relational/non-relational

## Powered by the cloud
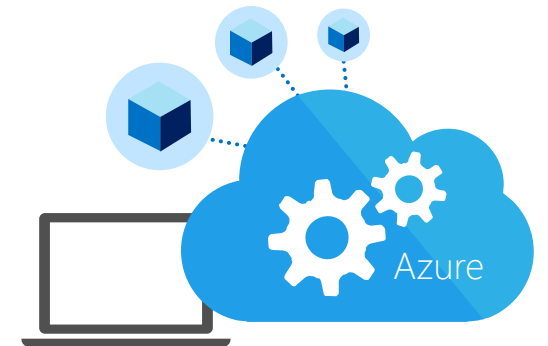
Get started in minutes

Integrate with Azure ML, Power BI, and ADF

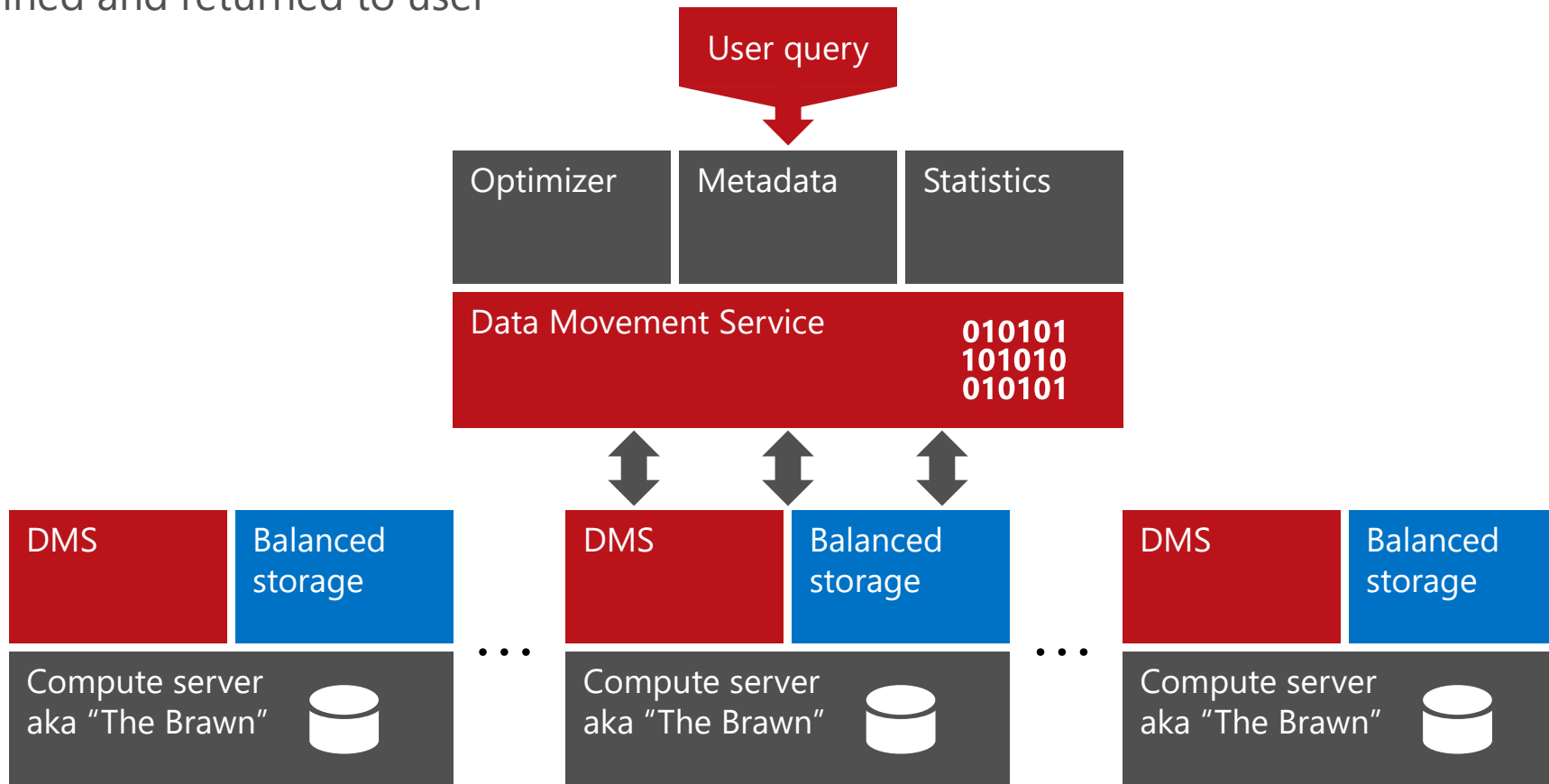Become enterprise-ready

SaaS

Azure

Public cloud

Office 365
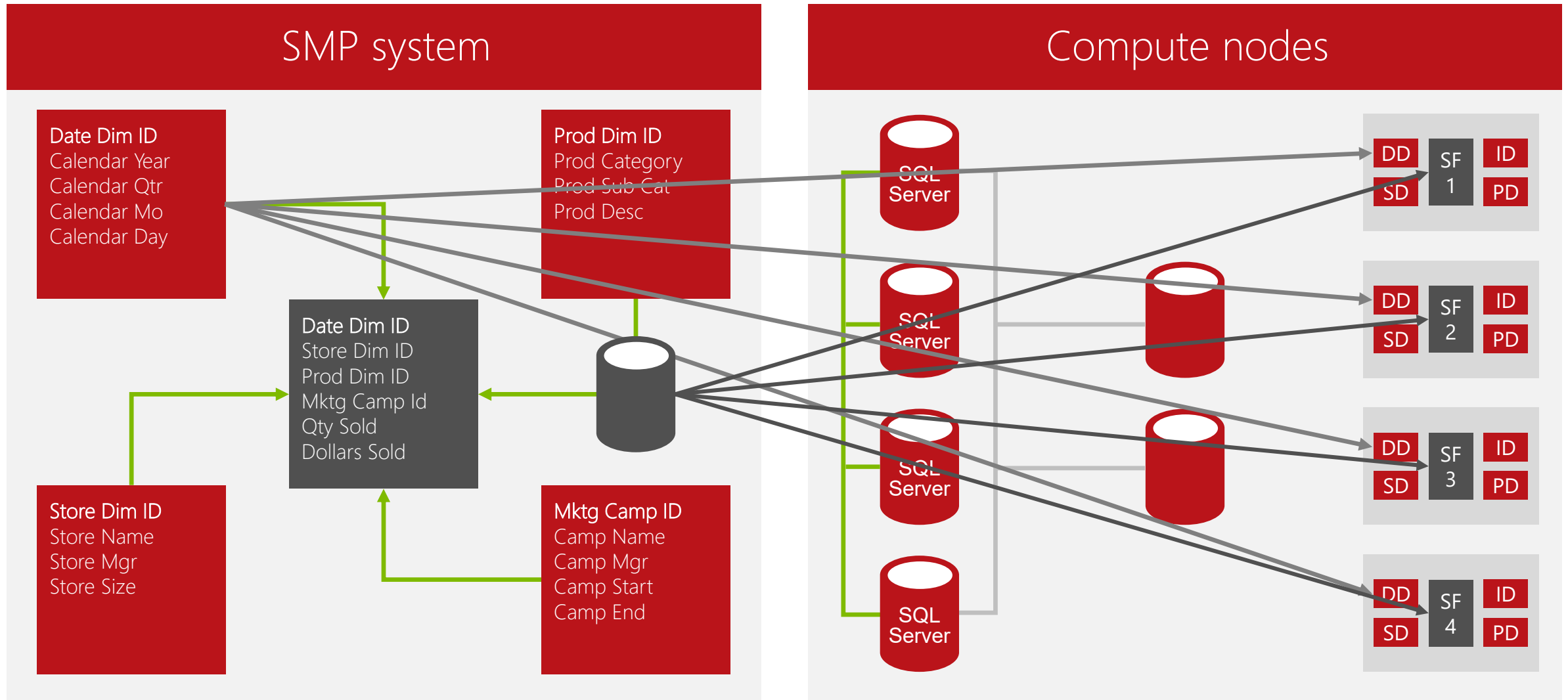
## Market-leading price & performance

Azure

Use simple billing compute and storage

Pay for what you need, when you need it with dynamic pause

Bring DW to the cloud without rewriting

# Logical architecture

1. Optimizer creates parallel query plan

2. Each compute server runs portion of query in parallel

3. Data is combined and returned to user

User query

| Optimizer | Metadata | Statistics |
|---|---|---|

Data Movement Service    **010101**
                         **101010**
                         **010101**

| DMS | Balanced storage |
|---|---|

Compute server aka "The Brawn"

...

| DMS | Balanced storage |
|---|---|

Compute server aka "The Brawn"

...

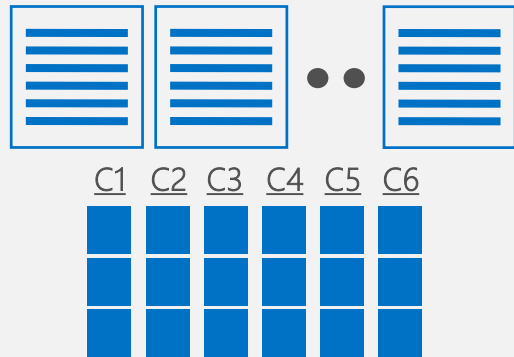| DMS | Balanced storage |
|---|---|

Compute server aka "The Brawn"

# MPP SQL table geometries

# Blazing-fast performance

## MPP and in-memory columnstore for next-generation performance

### Columnstore index representation



C1  C2  C3  C4  C5  C6

### Parallel query execution



Query

Results

Updateable clustered columnstore vs. table with customary indexing

Up to **100x** faster queries
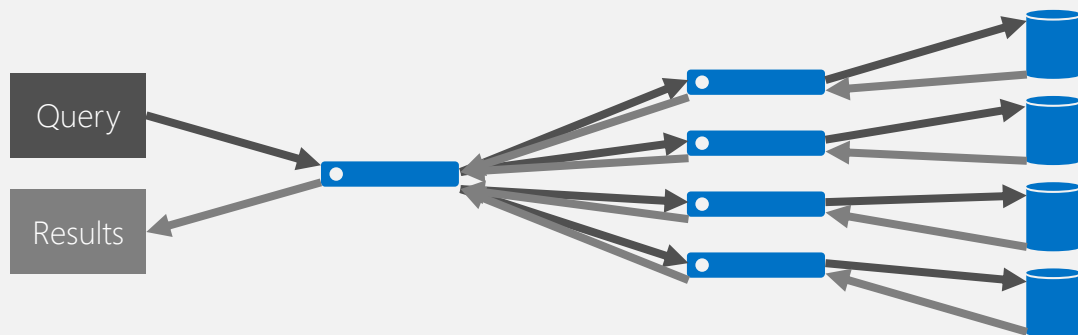
Up to **15x** more compression



Data storage in columnar format for massive compression

Data loading into or out of memory for next-generation performance, with up to 60% improvement in data loading speed

Updateable and clustered for real-time trickle loading

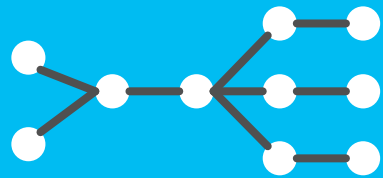# Better together: Azure SQL DW Service and APS

## Test/dev
Test new ideas in SQL Data Warehouse before rolling out to production in APS

## Age data
Age data to SQL Data Warehouse, but maintain full MPP power

## Company policy restrictions
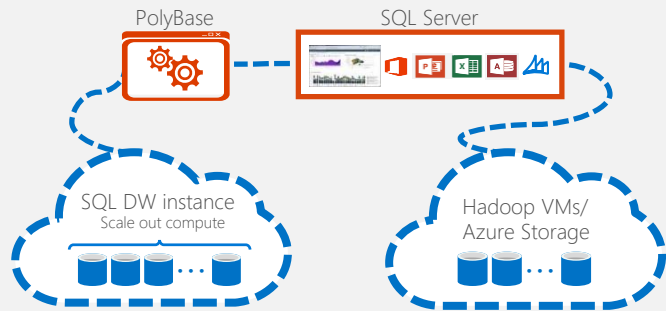Store data in APS that company policy prohibits from being in the cloud

## Disaster recovery
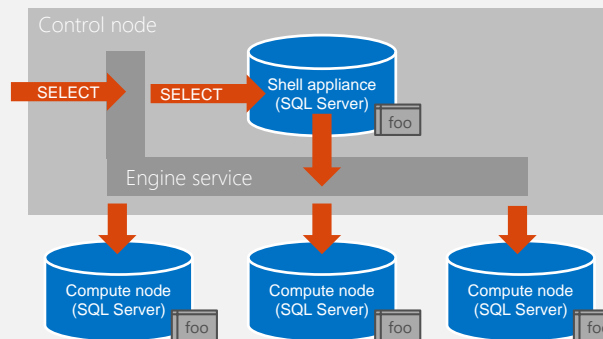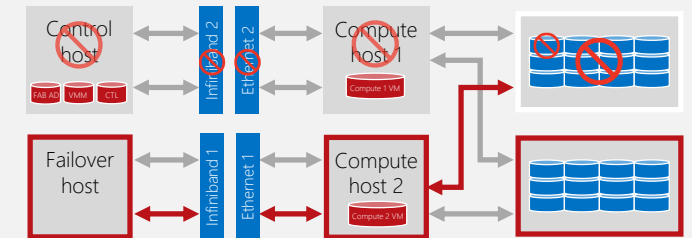Use SQL Data Warehouse or APS as disaster recovery solution with dual load

# Conclusion

## PolyBase

PolyBase      SQL Server

SQL DW instance
Scale out compute

Hadoop VMs/
Azure Storage

## Massively parallel processing

Control node

SELECT → SELECT → Shell appliance
(SQL Server)   foo

Engine service

Compute node
(SQL Server)   foo

Compute node
(SQL Server)   foo

Compute node
(SQL Server)   foo

## High availability

Control host
FAB AD   VMM   CTL

Infiniband 2   Ethernet 2

Compute host 1
Compute 1 VM

Failover host

Infiniband 1   Ethernet 1

Compute host 2
Compute 2 VM

---

## Microsoft APS

The Microsoft Analytics Platform System can meet the demands of your evolving data warehouse environment with its scale-out, massively parallel processing integrated system supporting hybrid data warehouse scenarios. It provides the ability to query across relational and non-relational data by leveraging Microsoft PolyBase and industry-leading big data technologies.

Azure SQL Data Warehouse enables APS customers with different workloads to leverage a cloud-based MPP engine and cloud-based analytics by supporting a hybrid architecture or eco-system with APS + Azure SQL Data Warehouse.

## Azure SQL Data Warehouse

# Big Data from Microsoft

# Azure HDInsight

A Cloud Spark and Hadoop service for the Enterprise

**Reliable** with an **industry leading SLA**

**Enterprise-grade security** and **monitoring**

**Productive platform** for **developers** and **scientists**

**Cost effective** cloud scale

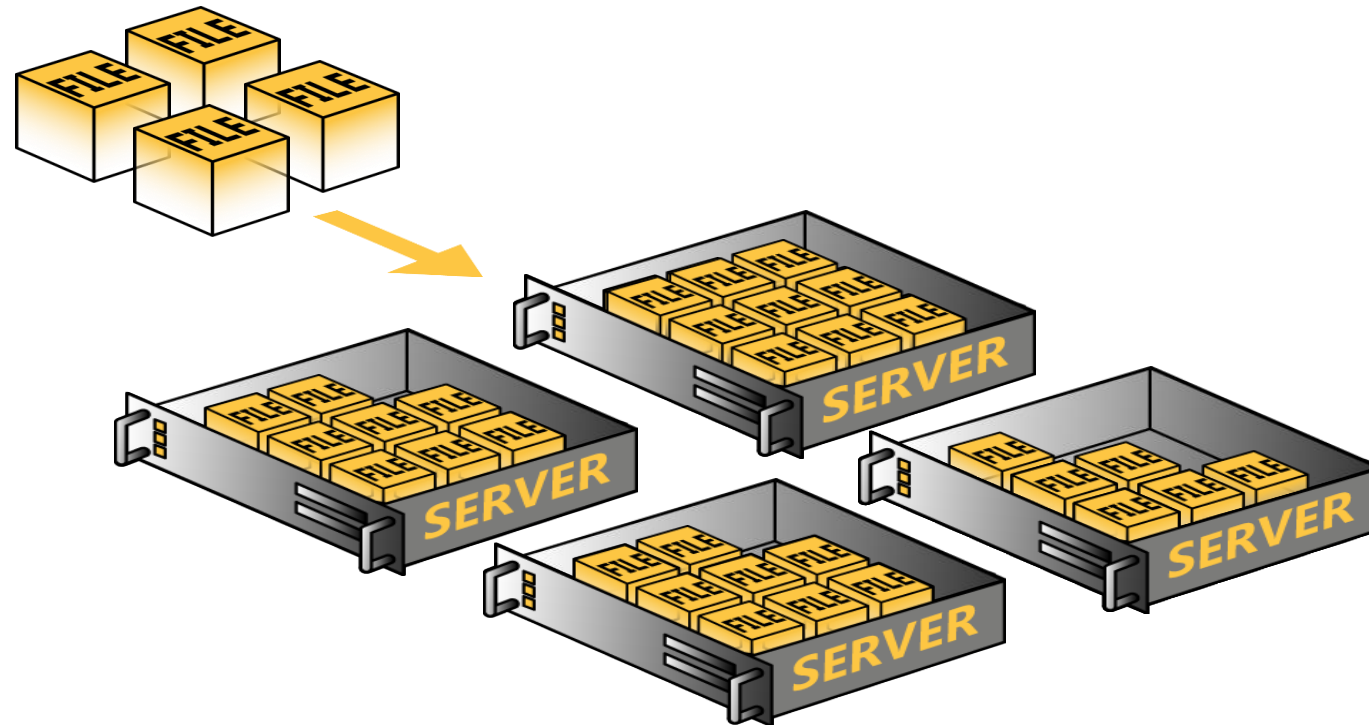**Integration** with leading **ISV applications**

**Easy** for administrators to **manage**

**63% lower TCO** than deploy your own Hadoop on-premises*

*IDC study "The Business Value and TCO Advantage of Apache Hadoop in the Cloud with Microsoft Azure HDInsight"

Microsoft

# So how does it work?
## First, store the data

# So how does it work?
## Second, take the processing to the data



```
// Map Reduce function in JavaScript

var map = function (key, value, context) {
var words = value.split(/[^a-zA-Z]/);
for (var i = 0; i < words.length; i++) {
            if (words[i] !== "")
{context.write(words[i].toLowerCase(), 1);}
}};

var reduce = function (key, values, context) {
var sum = 0;
while (values.hasNext()) {
sum += parseInt(values.next());
            }
context.write(key, sum);
};
```

# HDInsight Storage Infrastructure



Azure Blob Storage

Azure Flat Network Storage

HDInsight Compute Nodes
(Large VMs)

http://dennyglee.com/2013/03/18/why-use-blob-storage-with-hdinsight-on-azure/

# Recognized by top analysts

FORRESTER®



Challengers | Contenders | Strong Performers | Leaders

Strong

Microsoft

Amazon Web Services • IBM

Current offering

Google

Qubole • • Altiscale
Oracle

• Rackspace

Market presence

Weak

Weak —————— Strategy —————— Strong

## Forrester Wave for Big Data Hadoop Cloud

- Named industry leader by Forrester with the most comprehensive, scalable, and integrated platforms*

- Recognized for its cloud-first strategy that is paying off*

*The Forrester WaveTM: Big Data Hadoop Cloud Solutions, Q2 2016.

Microsoft

# Lower total cost of ownership

- No hardware

- Hadoop support included with Azure support

- Pay only for what you use

- Independently scale storage and compute

- No need to hire specialized operations team

- 63% lower total cost of ownership than on-premises*

*IDC study "The Business Value and TCO Advantage of Apache Hadoop in the Cloud with Microsoft Azure HDInsight"

# Azure
# Data Lake Store

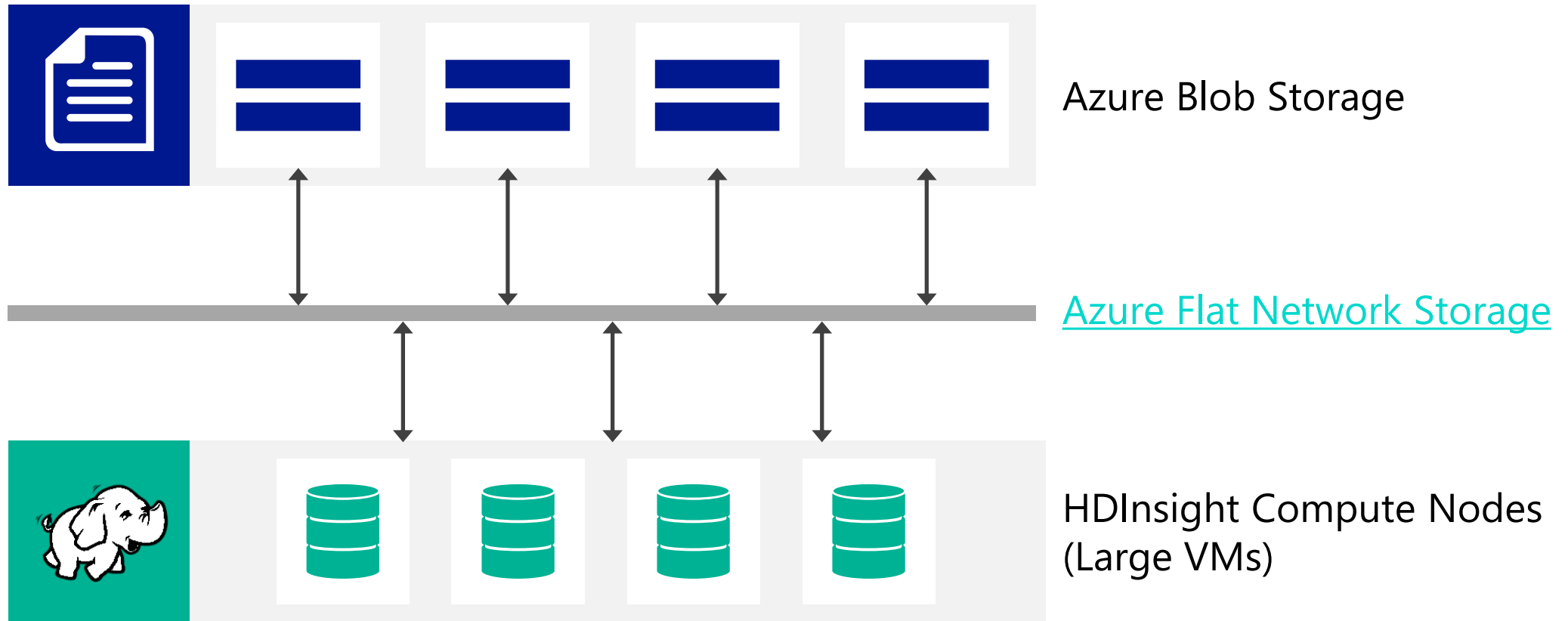A No limits Data Lake that powers Big Data Analytics

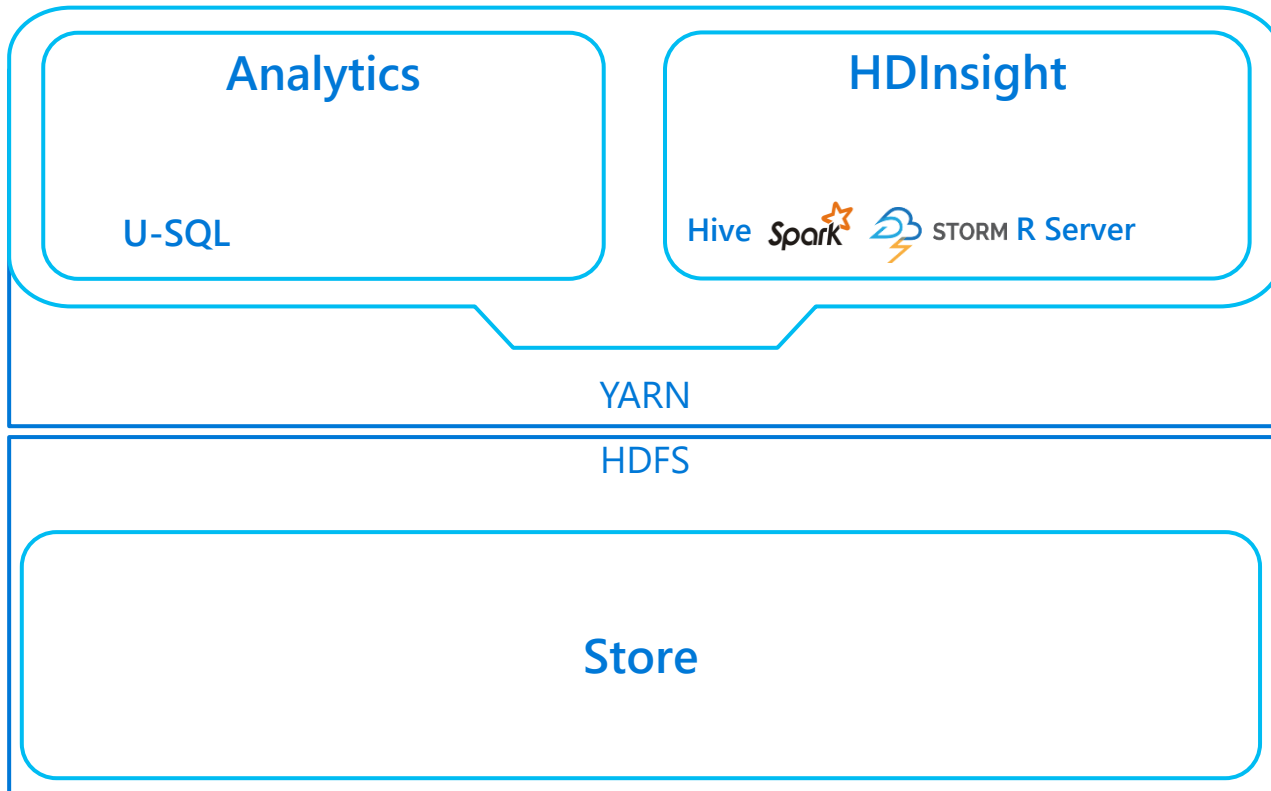**Petabyte size files** and **Trillions of objects**

**Scalable throughput** for **massively parallel analytics**

**HDFS** for the cloud

**Always encrypted, role-based security & auditing**

**Enterprise-grade** support

# Azure Data Lake

**Analytics**

U-SQL

**HDInsight**

Hive Spark STORM R Server

YARN

HDFS

**Store**

Store and analyze data of any kind and size

Develop faster, debug and optimize smarter

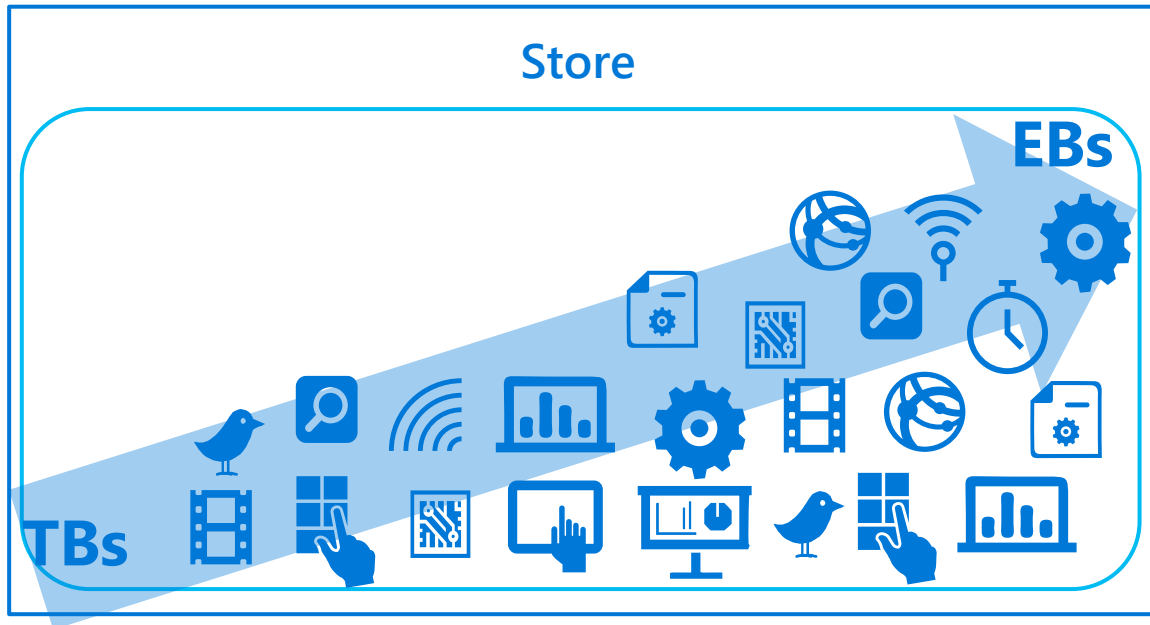Interactively explore patterns in your data

No learning curve

Managed and supported

Dynamically scales to match your business priorities

Enterprise-grade security

Built on YARN, designed for the cloud

# Petabyte size files and Trillions of objects



**Store**

EBs

TBs

- Store data in it's native format

- PB sized files, 200x larger than anyone else

- Scalable throughput  for massively parallel analytics

- No need to redesign application or reparation data at higher scale

Microsoft

# Anatomy of a U-SQL query

**10 log records by Duration (End time minus Start time). Sort rows in descending order of Duration.**

ClassLibrary2 - Microsoft Visual Studio

File  Edit  View  Project  Build  Debug  Team  SqlIP  Tools  Test  Analyze  Wi

Debug    Any CPU    ▶ Start

```
REFERENCE ASSEMBLY WebLogExtASM;

@rs =
    EXTRACT
        UserID          string,
        Start           DateTime,
        End             DatetTime,
        Region          string,
        SitesVisited    string,
        PagesVisited    string
    FROM "swebhdfs://Logs/WebLogRecords.txt"
    USING WebLogExtractor();

@result = SELECT UserID,
        (End.Subtract(Start)).TotalSeconds AS Duration
        FROM @rs  ORDER BY Duration DESC FETCH 10;

OUTPUT @result TO "swebhdfs://Logs/Results/top10.txt"
USING Outputter.Tsv();
```

**Rowset**: Conceptually is like an intermediate table… is how U-SQL passes data between statements

- U-SQL types are the same as C# types
- The structure (schema) is first imposed when the data is first extracted/read from the file (schema-on-read)

Input is read from this file in ADL

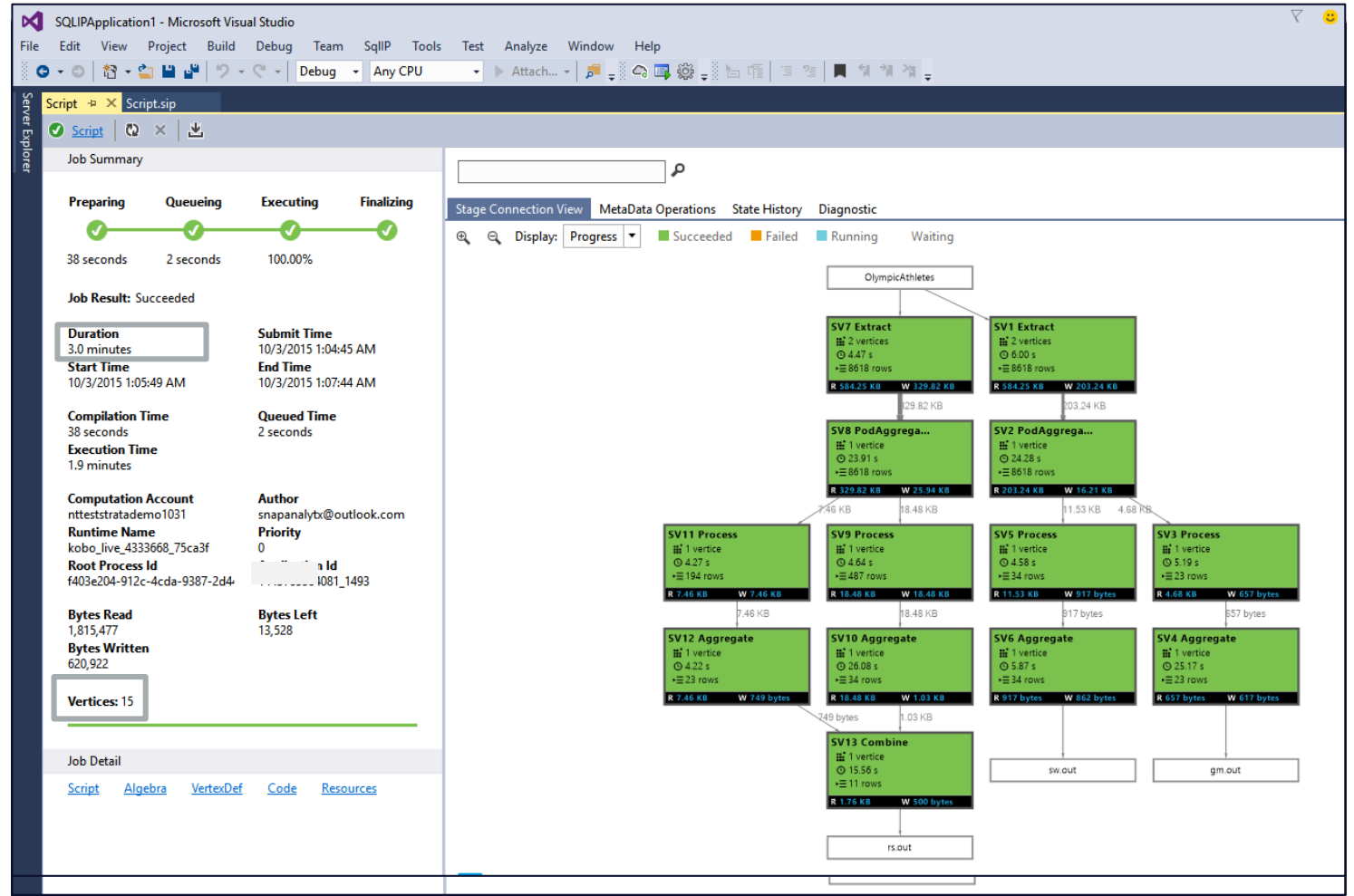Custom function to read from input file

C# Expression

Output is stored in this file in ADL

Built-in function that writes the output in TSV format

Microsoft

42

# Job execution graph

⚡ After a job is submitted the progress of the execution of the job as it goes through the different stages is shown and updated continuously

⚡ Important stats about the job are also displayed and updated continuously