**Les King**

**Director, Hybrid Data Management Solutions**

**May, 2018**

lking@ca.ibm.com

ca.linkedin.com/pub/les-king/10/a68/426

# Hybrid Data Management Strategy and New News !

# Les King

Director, Hybrid Data Management Solutions
Professor, Big Data, Data Warehousing and Db2, Seneca College

**lking@ca.ibm.com**
**ca.linkedin.com/pub/les-king/10/a68/426**

## Professional Highlights
- 27 years of Information Management, Database and Analytics
- Technical sales
- Technical customer support
- Software development
- Product / Offering management
- Product Marketing
- Product Sales
- Taught mathematics at University of Toronto
- Teaching data warehousing, big data and Db2 at Seneca College

# Safe Harbor Statement

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# All businesses have become data driven

# Data Professionals – Evolving Roles

*As data maturity increases, so does
the number of data professionals
who are hungry to put data to work*

**App Developers**
Easily plug into data and
models to make apps more
powerful

**Data Scientists**
Streamline algorithm
development to deliver
insight faster

**Business Analysts**
Easily discover and
explore data to
improve decisions

**Data Engineers**
Tame, curate and
secure data to
make it relevant
and accessible.

IoT

The Weather Company
An IBM Business

DBaaS

DB/DW

Public    Social

# The Challenges of Fast Data

**Data is arriving faster than ever before**
- Billions of events processed every day
- Evident cross industry and driven by IoT
- Must land data quickly, or throw it away

**Total data is large, and growing rapidly**
- Storing all events implies large data sets
- Storage costs are significant, and must be managed

**Data is useless without fast insights**
- Data value decays rapidly over time
- Insights must derived quickly, and use advanced analytics (ML)

**Data availability without duplication**
- Data must be available to the entire organization without requiring replication or duplication
- Maintain data in open format for future-proofing

AI

Machine Learning

Analytics

Data

**The "AI Ladder"**

## Data Management Strategy is HYBRID

**Its not about Cloud or On-Premises its about <span style="color:red">Cloud</span> <u>AND</u> <span style="color:red">On-Premises</span>**

**Its not about Traditional Relational or Open Source its about <span style="color:gold">Traditional Relational</span> <u>AND</u> <span style="color:gold">Open Source</span>**

<span style="color:#1f7fd0; font-size:2em">**It's About Hybrid**</span>

**Its not about SQL or NoSQL its about <span style="color:green">SQL</span> <u>AND</u> <span style="color:green">NoSQL</span>**

**Its not about Structured or Unstructured Data its about <span style="color:purple">Structured</span> <u>AND</u> <span style="color:purple">Unstructured</span> Data**

# Common SQL Engine – Business Value

A **COMMON SQL ENGINE** enabling true **HYBRID** data solutions for **ALL WORKLOAD** types

**Systems of Record**

**Systems of Engagement**

**Systems of Insight**

**Event Processing**

Investment Protection

**WRITE ONCE, RUN ANYWHERE**

**PUBLIC CLOUD**

**PRIVATE CLOUD**

**ON-PREMISES**

# Common SQL Engine – Consistent Technical Capabilities

A **COMMON SQL ENGINE** enabling true **HYBRID** data solutions for **ALL WORKLOAD** types

| **NEW** | | | **NEW** | **NEW** | | |
|---|---|---|---|---|---|---|
| Event Store | Managed Public Cloud Service | Hosted Public Cloud Service | On-premises Private Cloud | On-premises Appliance | On-premises Custom Software | Hadoop / Spark Environment |
| Db2 Event Store | Db2 [Warehouse] On Cloud | Db2 Hosted | Db2 Warehouse Db2 OLTP | IBM Integrated Analytics System | Db2 | Big SQL |

## Foundation

- ✓ Full MPP scalability (GB-PB)
- ✓ High Concurrency
- ✓ Load and Go Simplicity
- ✓ Consistent Management and WLM
- ✓ HA, DR & Replication
- ✓ Integrated Security & Encryption

## Application

- ✓ Built-in analytics (OLAP)
- ✓ Data Virtualization
- ✓ Application portability
- ✓ Hybrid by design
- ✓ Oracle Compatibility
- ✓ Netezza Compatibility

## New Growth Trends

- ✓ Spark Integration
- ✓ HTAP Support
- ✓ SQL & NOSQL Capabilities
- ✓ Native JSON Support
- ✓ R Language Support
- ✓ Structured & Unstructured Data

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# DB2 - Highlights and Strategic Investment Areas

**Petabyte Scale In-Memory Warehousing - Systems of Insight**

Build or expand your operational warehouse with petabyte scale BLU Technology

**Best Database for SAP**

Leverage the highest performing, most optimized and most cost effective solution for SAP

**Alternative to Oracle**

Minimized cost and risk of migration when breaking free from high costs of Oracle

**SQL and NOSQL Integration**

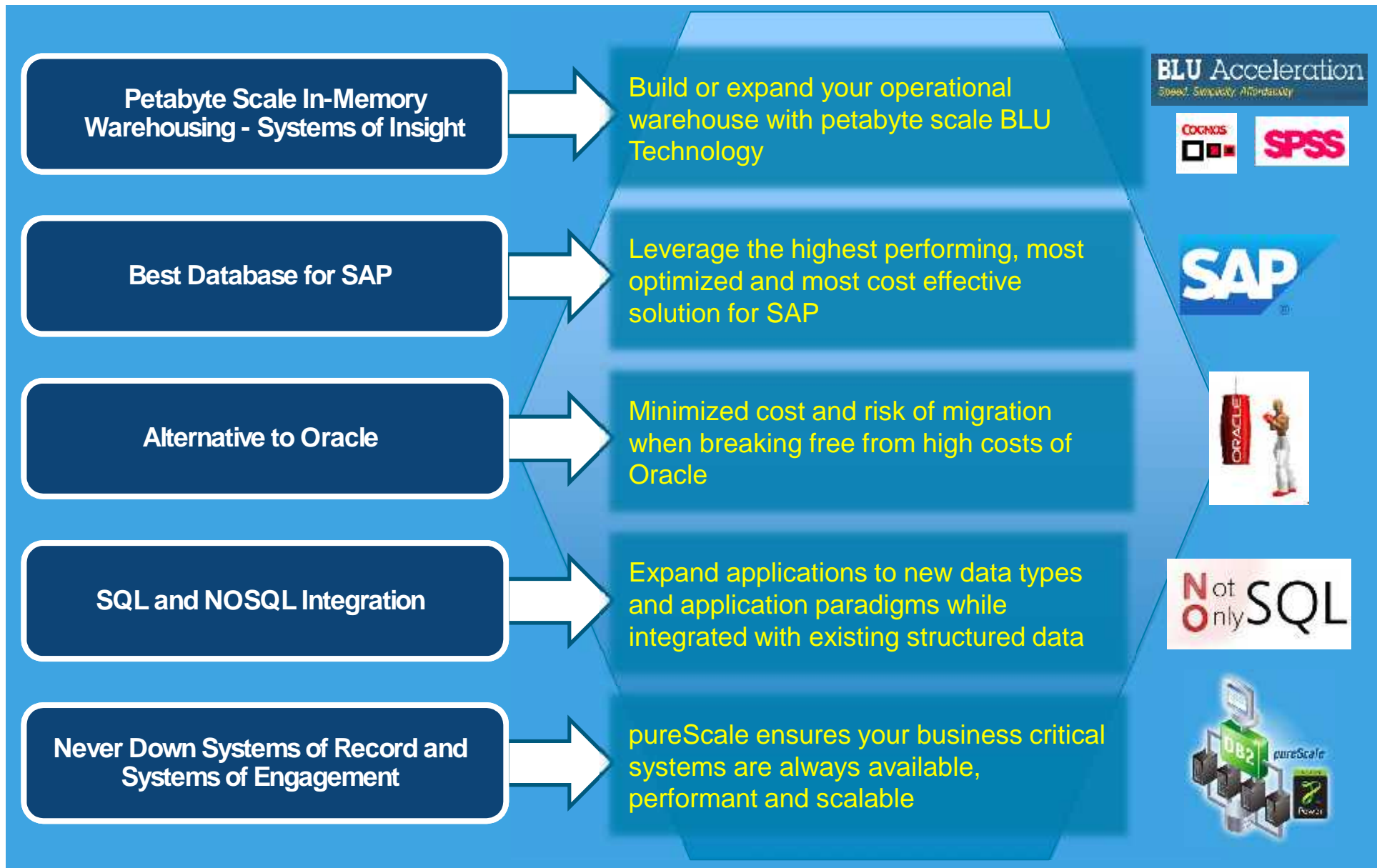Expand applications to new data types and application paradigms while integrated with existing structured data

**Never Down Systems of Record and Systems of Engagement**

pureScale ensures your business critical systems are always available, performant and scalable

**BLU** Acceleration
*Speed. Simplicity. Affordability.*

COGNOS   **SPSS**

**SAP**

ORACLE

**N**ot **O**nly SQL

DB2 *pureScale*

**Check George Baklarz's Presentation this afternoon**

# Db2 Version 11.1.2.2 Highlights

## Higher Availability and Core Capabilities

**Near-zero outage recovery**
- Online crash recovery
- pureScale REBUILD restore

## Column-Organized (BLU) Tables

**Deeper BLU Optimizations for Operational Workloads**
- Performance enhancements
- Builds on 4Q '16 advances
- Enables use of BLU beyond strictly analytic workloads

## NoSQL Support

**Native JSON support**
- JSON SQL support Part 1
- Built-in UDFs for enhanced JSON capabilities

N<sup>ot</sup> O<sup>nly</sup>SQL

## Additional Operating System Support

Solaris Support – by exception

MacOS Support – by exception

## Db2 Tooling Capabilities

- Data Server Manager
- DB2 Connect
- DS Driver
- DS Gateway
- Advanced Recovery Tools

## Packaging Changes

- Developer Community Edition
- Introduction of non-production licenses
- Data Management Bundle V1

Check George Baklarz's Presentation this afternoon

# Db2 Version 11.1.3.3 Highlights

## Higher Availability and Core Capabilities

- Faster Rollback of very large transactions
- WLM – Improve deadlock detection
- HADR Resilience and SSL Encryption
- Db2iupdt – ADD/DROP CFs on-line
- pureScale – on-line CREATE INDEX w/R/W access to table
- pureScale – faster member crash recovery

## Column-Organized (BLU) Tables

**UDF Cacheing for BLU**
BLU Memory Usage enhancements
Temporal Query Support
Index Support

## Additional Operating System Support

Solaris Support – 11.3+

## Data Virtualization

MariaDB Connectivity Support
Db2 iSeries 7.2&7.3 Connectivity Support
Teradata 16 Connectivity Support
JSON over RESTful Service (MongoDB)
Boolean, Binary/Varbinary Data Type Mapping Enhancement
Pushdown Improvement for Hadoop Datasource
Function Mapping Pushdown Enhancement
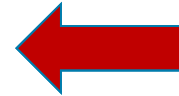
## Packaging Changes

- Hybrid Data Management Packaging

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# Db2 and the Cloud

**Provisioning & Db2 Setup**    **Management**

**Maintenance**

## "Bring Your Own License"
- Custom-deployable software on your own infrastructure or private cloud or public cloud
- Fully customizable for any type of workload
- Complete flexibility including DPF and pureScale *
- Customer managed

## Db2 Hosted
- Hosted database-as-a-service
- Pre-defined hardware configurations
- Fully customizable for any type of workload
- Available on SoftLayer and AWS
- Customer managed

## Db2 on Cloud
- Fully managed database-as-a-service
- Pre-defined and flexible hardware configurations optimized for transactional and general purpose workloads
- Available on Bluemix public cloud

## Db2 Warehouse on Cloud
- Fully managed database-as-a-service
- Pre-defined hardware configurations optimized for analytics workloads
- In-database analytics
- Available on Bluemix and AWS public cloud

## Db2 Warehouse
- Deploy on your own infrastructure or private cloud
- Docker container technology for fast and simple deployment
- Optimized for analytic workloads
- Scalable, elastic
- Customer managed

## Db2 OLTP
- Deploy on your own infrastructure or private cloud
- Docker container technology for fast and simple deployment
- Optimized for operational and OLTP workloads
- Scalable, elastic
- Customer managed

# Db2 and the Cloud

| | Provisioning & Db2 Setup | Maintenance | Management |
|---|---|---|---|

### "Bring Your Own License"

- Custom-deployable software on your own infrastructure or private cloud or public cloud
- Fully customizable for any type of workload
- Complete flexibility including DPF and pureScale *
- Customer managed

### Db2 Hosted

- Hosted database-as-a-service
- Pre-defined hardware configurations
- Fully customizable for any type of workload
- Available on SoftLayer and AWS
- Customer managed

### Db2 on Cloud

- Fully managed database-as-a-service
- Pre-defined and flexible hardware configurations optimized for transactional and general purpose workloads
- Available on Bluemix public cloud

### Db2 Warehouse on Cloud

- Fully managed database-as-a-service
- Pre-defined hardware configurations optimized for analytics workloads
- In-database analytics
- Available on Bluemix and AWS public cloud

### Db2 Warehouse

- Deploy on your own infrastructure or private cloud
- Docker container technology for fast and simple deployment
- Optimized for analytic workloads
- Scalable, elastic
- Customer managed

### Db2 OLTP

- Deploy on your own infrastructure or private cloud
- Docker container technology for fast and simple deployment
- Optimized for operational and OLTP workloads
- Scalable, elastic
- Customer managed

# Introducing IBM Cloud Private

| Innovation | Integration | Investment Protection | Management and Compliance |
|---|---|---|---|

Kubernetes-based container platform

Cloud Foundry for prescribed container-based application development and deployment and life cycle management

Integrated DevOps toolchain

Catalog of integration services

API availability and management to integrate applications and data across environments

Prescriptive guidance on where to run and how to architect your critical workloads

Next generation versions of industry leading IBM Middleware and Analytics
(MQ, Db2, Data Science, Cognos, Blockchain, IIB)

Core operational services, including monitoring, log mgmt, and security

Integration with existing systems and operations management solutions

# Analytics Roadmap : Offerings / Capabilities on ICp

**( as of Nov 2017)**

Preliminary & Subject to changes
* To be confirmed

## 2017 Q4

- Db2 OLTP

- Db2 Warehouse

- Data Server Manager

- Data Science Experience

## 2018 1H

**1. Hybrid Data Management**

- Db2 OLTP

- Db2 Warehouse MPP

- Data Server Manager

- Big SQL *

**2. Unified Governance**
- Data Stage
- IGC

**3. Data Science & BA**
- Data Science Experience

**Common / Foundational**

| ✓ | Metering | ✓ | Logging |
|---|---|---|---|
| ✓ | Monitoring | ✓ | IAAM & SSO |
| ✓ | Catalog | | |

## 2018 2H

**1. Hybrid Data Management**

- Db2 OLTP MPP

- Db2 Event Store

**2. Unified Governance**

- WEX *

**3. Data Science & BA**
- SPSS Modeler *
- SPSS Statistics *
- Cognos *

**Common / Foundational**

| ✓ | Metering | ✓ | Logging |
|---|---|---|---|
| ✓ | Monitoring | ✓ | IAAM & SSO |
| ✓ | Catalog | | |

# Why Analytics on IBM Cloud Private

**True Hybrid Solution** - consistency between public cloud and private cloud

No vendor lock-in. **Open Platform as a Service** (PaaS) for maximum integration ability

**Container-based platform** with very fast time to value (hours instead of weeks)

Extensive service-oriented analytic and machine learning capabilities ready for **Data Scientists** and **Business Analysts**

Optimized and secure **Data Management Services** for SQL, NoSQL, structured, semi-structured and unstructured data

**Secure, governed** and **compliant** platform for integration with any data source

# IBM Cloud Private – More Information

Catch Kelly Schlamb's session this afternoon !!

**Learn More**

- IBM Cloud private home: https://ibm.biz/Bdj4Bz

- White paper: https://ibm.biz/Bdj4UJ

**See it In Action**

- Offering demo: https://youtu.be/yzXA3qhfaq0

- Try It: https://ibm.biz/Bdj4UC

- Free Community Edition: https://hub.docker.com/r/ibmcom/cfc-installer/

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# Portfolio Simplification:
## *Three new bundles*

### Hybrid Data Management

- Db2
- Db2 Warehouse
- Db2 Event Store
- Db2 Big SQL

### Unified Governance & Integration

- Information Server
- Entity Analytics
- Master Data Management
- Info Governance Catalog
- Data Replication
- Test Data Fabrication
- Info Lifecycle Governance
- Industry Models
- BigIntegrate, BigQuality BigMatch

### Data Science & Business Analytics

- DSx Local
- SPSS Modeler
- Decision Optimization
- Watson Explorer (v12 +)
- Cognos Analytics
- Planning Analytics

**We will now focus on Hybrid Data Management**

# FlexPoints: How It Works

**Buy FlexPoint licenses for the "Platform of your Choice"**



Flex Points

Choose how you want to deploy your FlexPoints

Each component has a FlexPoint price

Db2      Big SQL      Db2 WH

Swap components as your needs change

Platform Offerings deliver integrated capabilities – now offered as flex bundles to simplify planning for adoption and growth at the lowest cost

**Available for Our 3 Platform Offerings:**

➢ Hybrid Data Management

  ➢ Db2

  ➢ Db2 Warehouse

  ➢ Db2 Event Store

  ➢ Db2 Big SQL

➢ Unified Governance & Integration

➢ Data Science & Business Analytics

**FlexPoints CANNOT be used across PLATFORMS**
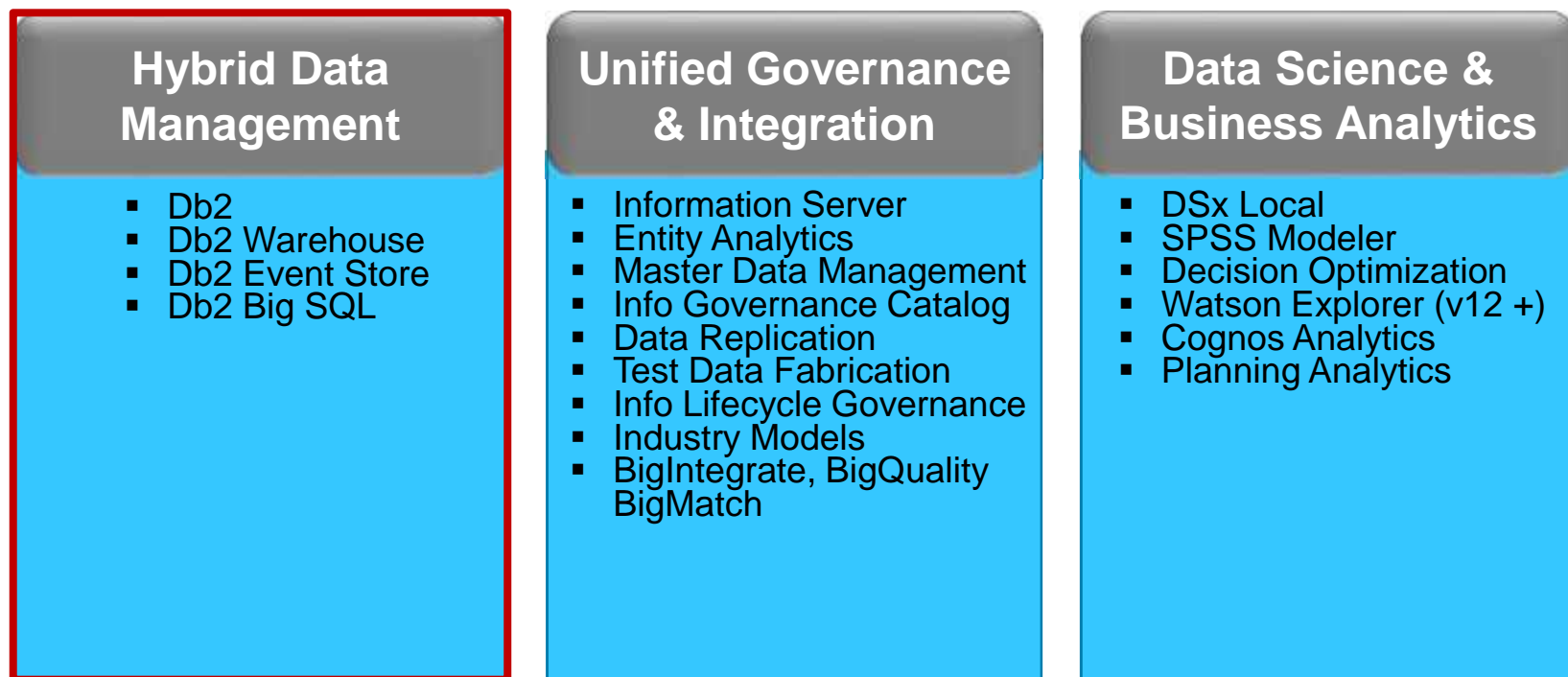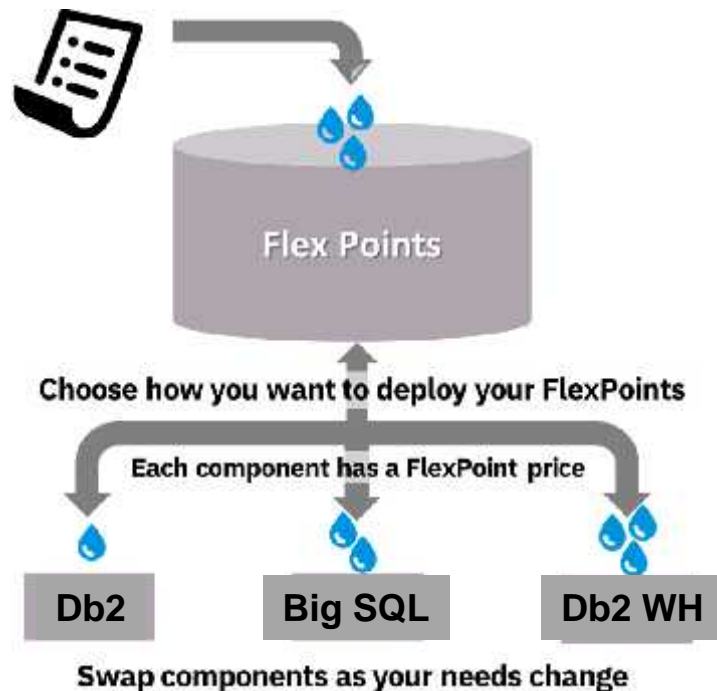As an example, Data Science and Business Analytics FlexPoints are NOT valid for Hybrid Data Management

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# Next Generation Analytics Appliance

World's fastest and "greenest" analytical platform

World's First Analytic Data Warehouse Appliance

World's First Petabyte Data Warehouse Appliance

World's First 100 TB Data Warehouse Appliance

World's First Data Warehouse Appliance

PureData System for Analytics N3000

PureData System for Analytics N2000

TwinFin™ with i-Class™ Advanced Analytics

TwinFin™

NPS® 10000 Series

NPS® 8000 Series

27

| 2003 | 2006 | 2009 | 2010 | 2012 | 2014 | 2017+ |

# Next Generation Analytics Appliance – Names

**4000-series**
**"Sailfish"**
**IBM Integrated Analytics System**
**IIAS**

World's fastest and "greenest" analytical platform

World's First Analytic Data Warehouse Appliance

World's First Petabyte Data Warehouse Appliance

World's First 100 TB Data Warehouse Appliance

World's First Data Warehouse Appliance

New

PureData System for Analytics N3000

## Our Next-Generation Hybrid DWH vision:

**Sailfish** is IBM's industry-leading hybrid - Private Cloud Data Warehouse and Analytics Platform that will **integrate** seamlessly with other ground and cloud data warehouse services, delivering ultra **fast & scalable** performance, cloud **elasticity** together with end to end security - and the ultimate in **simplicity** across all dimensions of the client's experience.

...Integrated, Insightful, Simple, Agile and Secure!

28

2003    2006    2009    2010    2012    2014    2017+

# IBM Integrated Analytics System

## *Next Generation Hybrid Data Warehouse*

Optimized for **high performance** to support the broadest array of workload options for structured and unstructured data in your **hybrid data management** infrastructures

**Reliable, elastic and flexible** system that reduces and **simplifies management** resources

Real time analytics with **machine learning** that accelerates decision making, bringing new opportunities to the business – ready for **business analysts** and **data scientists**

Leverages a **Common SQL Engine** for workload portability and skill sharing across public and private cloud

**Cloud-ready** to support multiple workload deployment options

Built-in **IBM Data Science Experience** to collaboratively analyze data

# Addressing Top Customer Requirements

**Broader set of workloads**
- Combination of reporting, analytics, operational analytics and data stores

**Higher Concurrency**
- Expand number of business analytics and machine learning activities within a single system

**In-Place Expansion**
- Independently scale both compute and storage as needed while protecting existing investments

**Richer Availability Solutions**
- High Availability, Disaster Recovery and replication solutions

# Less admin & more analytics

Simplicity

**Accelerate Time to Insight**
Easy to Deploy and Easy to Operate
Faster Time to Value - Load and Go…it's an appliance!
Lower Total Cost of Ownership
Built-in Tools for data migration and data movement

**Load and Go**

**BI Developers & DBAs – faster delivery times**
No configuration
No storage administrations
No physical modeling
No indexes and tuning
Data model agnostic
Self Service Management dashboard

**Low TCO**

**One Touch Support**

**ETL Developers**
No aggregate tables needed – simpler ETL logic
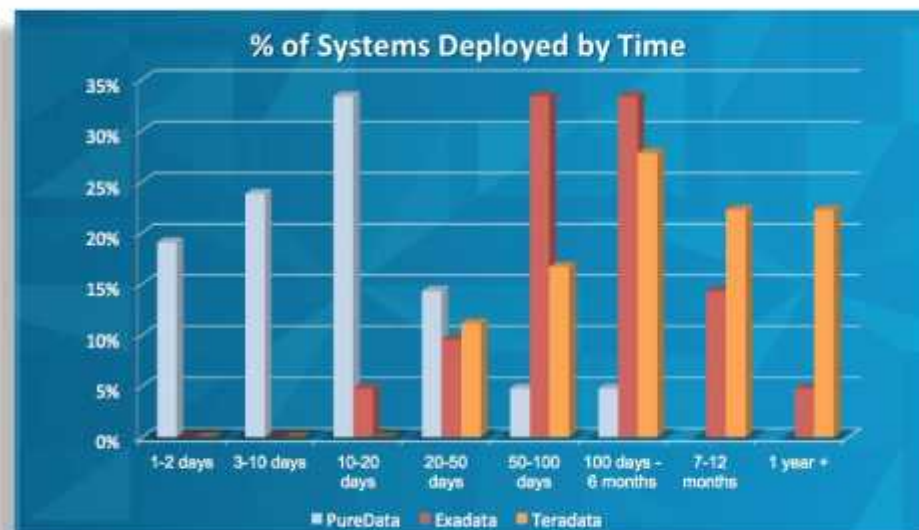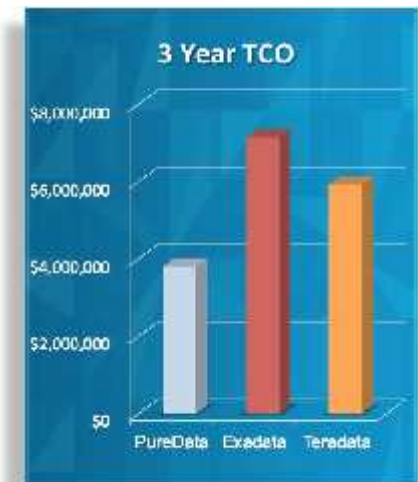Faster load and transformation times

**Business Analysts**
True ad hoc queries – no tuning, no indexes
Ask complex queries against large datasets
Load & query simultaneously

# Maintain Core Values



-Reduced administration
-Performance portal
-Lower end starting point
-More scale-out
-Fast time to deployment
-Low TCO

# Speed of Thought Analytics

**Performance**

**2X – 5X
Performance
Gain**

POWER

**Powered by RedHat® Linux on Power**

Optimized for Analytics with 4X Threads per core, 4X Memory bandwidth and 4X more cache at lower latency compared to x86

**ALL Flash Storage**

Hardware Accelerated architecture enabling faster insights with extreme performance, 99.999% reliability and operational efficiency

**MPP Scale out**

**Memory Optimized**

In-memory BLU columnar processing with dynamic movement of
data from storage

**Data Skipping**

Skips unnecessary processing of irrelevant data

**Actionable Compression**

Patented compression technique that preserves order so data can be used without decompressing

# Optimized Analytics Performance

## Next Generation In-Memory

In-memory columnar processing with dynamic movement of data from storage

## Analyze Compressed Data

Patented compression technique that preserves order so data can be used without decompressing

Encoded

## Embedded Spark

Spark As an Analytics Engine

Spark/R, Spark/ML, Rest API, Object Store ETL, Complex Transformations (ELT), Streaming

## CPU Acceleration

Multi-core and SIMD parallelism (Single Instruction Multiple Data)

**Instructions** **Data**

**Results**

## Data Skipping

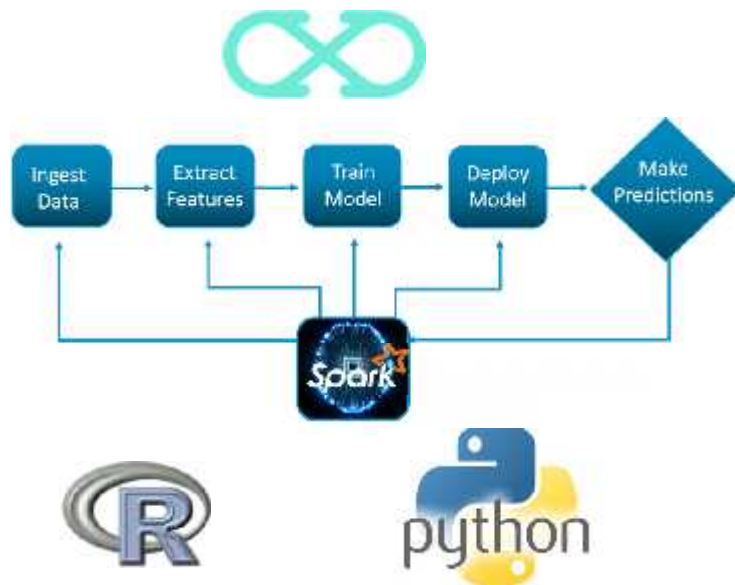Skips unnecessary processing of irrelevant data

## Powered by Hardware

Designed for Deep Complex Analytics

4X Threads per core
4X Memory Bandwidth
4X More cache at Lower Latency

# Ready for Data Scientists and Business Analysts

## Machine Learning

**Integrated Cognitive Assist for Machine Learning**
**DSX for Interactive & Collaborative Data Science**
Scalable ML Model Training, Deployment and Scoring with
Spark embed Predictive / Prescriptive In place Analytics

**Embedded**

Data mining, prediction, transformations, statistics, geospatial,
data preparation

**Full integration with tools for BI & visualization**

IBM Cognos, Tableau, Microstrategy, Business Objects, SAS,
MS Excel, SSRS, Kognitio, Qlikview

**Full integration with tools for model building and scoring**
IBM SPSS, SAS, Open Source R, Fuzzy Logix

**Full integration for custom analytics**

Open Source R, Java, C, C++, Python, LUA

# DSX Local on IIAS Benefits

- **The inclusion of DSX Local widens the audience for IIAS**
  - DSX Local is a on-prem platform which manages and provides access to the data, tools and packages that data scientist need
    - Jupyter, Zeppelin*, and RStudio
    - Anaconda for Python 2 and 3* support
    - Support for Python, Scala, and R languages

- **DSX Local extends IIAS federation support**
  - Livy included for connecting to and running jobs on external Spark clusters
  - GUI for connecting to external data sources and data sets
    - DB2, DB2 Z, Netezza, Informix, Oracle, dashDB, HDFS*, Hive*, and more to come
  - Easily combines data from multiple sources to create new data sets

- **DSX Local provides full model management for IIAS**
  - Create models with the built-in model builder GUI or programmatically from a notebook

# Write Once, Run Anywhere

## Hybrid



IBM Data Lift



IBM Fluid Query

### Application Agility

Common SQL Engine with comprehensive tools and capabilities across all deployment models: Public/Private Cloud, On-premise Appliance.

One ISV certification for all deployments .

### Operational Compatibility

Single consistent interface powered by IBM Data Server Manager for Management and Maintenance

### Make Data Simple and Accessible to All

Data Virtualization capabilities enabled by Fluid across deployment models

Querable Archive  Query historical data on Hadoop or other content stores

Discovery & Exploration Implement the Logical Data Warehouse; Land data in Hadoop for discovery, exploration & "day 0" archive

Build Bridges to RDBMS Islands Combine data from different enterprise divisions currently trapped in silos ; Federate to other data sources such as Oracle, SQL Server, PostgreSQL, Teradata, etc.,

### Ground to Cloud Blazing-fast Data Transfer

Integrated high speed IBM Data Lift using IBM Aspera for secure ground to cloud data movement

# Hybrid – Common SQL Engine

# Hybrid – Common SQL Engine


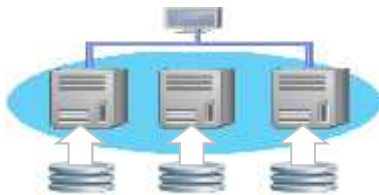
Db2 Warehouse
On Cloud

Db2 Big SQL

Db2 Warehouse

IBM Integrated
Analytics System

Db2 Big SQL

# Unmatched multi-dimensional Flexibility



**Flexible**

Scalable

Versatile Workloads

### In-Place Incremental Expansion
Easily and incrementally scale out your environment by adding Compute and Storage capacity to meet your growth needs

### In-place Tiered Storage Expansion
Independently scale storage for cost effective capacity growth

### HTAP with IBM Db2 Analytics Accelerator
Seamlessly integrate with IBM z Systems infrastructure to enable real-time analytics combining transactional data, historical data and predictive analytics

### Truly a Mixed Workload Appliance
Whether it be high scan performance needed to answer your business's strategic questions, high concurrency, low-latency requirements to support your operational systems, or even use as an operational data store. Perform all your enterprise Analytics needs on a single platform with mission critical availability.

### Flexible Licensing
Flexible entitlements for business agility & cost-optimization

# Expansion capabilities

**Non-disruptive in-place incremental expansion**

- **Reduce disruptions to your analytics systems as you scale out**

**Cloud-ready**

- **Tools to move workloads seamlessly to the cloud based on your requirements**

**Non-disruptive in-place tiered storage expansion**

- **Independently scale storage for cost effective capacity growth**

**Cost efficient multi-temperature storage**

- **Most frequently accessed data ("hot") on faster flash storage**
- **Less frequently accessed data ("colder") on cost efficient enterprise storage systems**

# IBM Db2 Analytics Accelerator

## High performance for complex queries

- Unprecedented response times to enable 'train of thought' analyzes frequently blocked by poor query performance

## Seamless integration with z Applications

- Brings high performance queries to existing z systems while protecting the core OLTP workloads

## Self-managed workloads

- Queries are executed in the most efficient location

## Transparent application access

- Brings the value of the Common SQL Engine to the z environment
- Applications connected to Db2 are entirely unaware of the Accelerator, all security is handled by Db2 z/OS

## Fast deployment and time to value

- Non-disruptive installation. Plug it in, load data and go in 1-2 days
- Db2 for z/OS query router automatically sends analytic queries to source which will provide optimal performance

**IBM Db2**

**A high performance appliance that integrates the IBM Integrated Analytics System with zEnterprise technology to deliver dramatically faster business analysis**

# One API – One implementation – Two deployment options

**Hardware Appliance**

**Deployment on IBM Z**



Uniform experience, simultaneous use, and easy transition between different implementations

Common analytics engine across all the platforms: Db2 Warehouse

# IBM Integrated Analytics System configurations

IBM Power 8 S822L  24 core server 3.02GHz
IBM FlashSystem 900

In-place Expansion Tiered storage

Mellanox 10G Ethernet switches
Brocade SAN switches

| | M4001-003 1/3 Rack | M4001-006 2/3 Rack | M4001-010 Full Rack | M4001-020 2 Racks | M4001-040 **4 Racks** |
|---|---|---|---|---|---|
| Servers | 3 | 5 | 7 | 14 | 28 |
| Cores | 72 | 120 | 168 | 336 | 672 |
| Memory | 1.5 TB | 2.5 TB | 3.5 TB | 7 TB | 14 TB |
| User capacity (Assumes 4x compression) | 64 TB | 128 TB | 192 TB | 384 | 768 |
| Tiered storage (Optional) | TBD—GA 1H 2018 | | | | |

**2 Racks + Tiered Storage targeted for 1H 2018;  In place expansion targeted for 2H 2018**

# Hardware architecture overview

**7 Compute Nodes in 1 rack containing**
- IBM Power 8 S822L  24 core server 3.02GHz
- 512 GB of RAM (each node)
- 2x 600GB SAS HDD
- Red Hat® Linux OS

**Up to 3 Flash Arrays in 1 rack containing**
- IBM FlashSystem 900
- Dual Flash controllers
- Micro Latency Flash modules
- 2-Dimensional RAID5 and hot swappable spares for high availability

**2x Mellanox 10G Ethernet switches**
- 48x10G ports
- 12x40/50G ports
- Dual switches form resilient network

**IBM SAN64B 32G Fibre Channel SAN**
- 16Gb FC Switch
- 48x 32Gb/s SFP+ ports

User Data Capacity:
**192 TB***
(Assumes 4x compression)

Power Requirements:
**9.4 kW**

Cooling Requirements:
**32,000 BTU/hr**

Scales from:
**1/3rd Rack to 8 Racks**
(initial GA is 1/3rd to 1 Rack)

# Application and Operational Compatibility

## … compared to Netezza

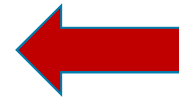| Perspective | September, 2017 (First release of Sailfish) | December, 2018 (Completion of Sailfish) |
|---|---|---|
| Applications | 95% SQL compatibility<br>nz* commands not available yet<br>Manual conversion of stored procedures<br>Performance degradation of INZA functions | 100% Application Compatibility<br>Equal or better performance for all applications |
| Operations & Management | nz* commands not available yet<br>Workload Management tools change<br>Replication solutions changed (NRS)<br>Multi-tenancy (single database) | Some areas of operational management will continue to be different with Sailfish in order to provide a richer set of capabilities (WLM, HA/DR) |

## … compared to PDOA and ISAS (Db2)

| Perspective | September, 2017 (First release of Sailfish) | December, 2018 (Completion of Sailfish) |
|---|---|---|
| Applications | 100% compatibility | 100% compatibility |
| Operations & Management | Some limitations such as multi-tenancy | 100% compatibility |

## … compared to Oracle

| Perspective | September, 2017 (First release of Sailfish) | December, 2018 (Completion of Sailfish) |
|---|---|---|
| Applications | 95%-98% compatibility<br>Leverages Oracle Application Compatibility Layer | 95%-98% compatibility<br>Leverages Oracle Application Compatibility Layer |

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# IBM and Hortonworks Deliver Data Science at Scale

*Focus on extending data science and machine learning to analyze the data in Apache Hadoop systems*

*Consumers get the best in class open technology*
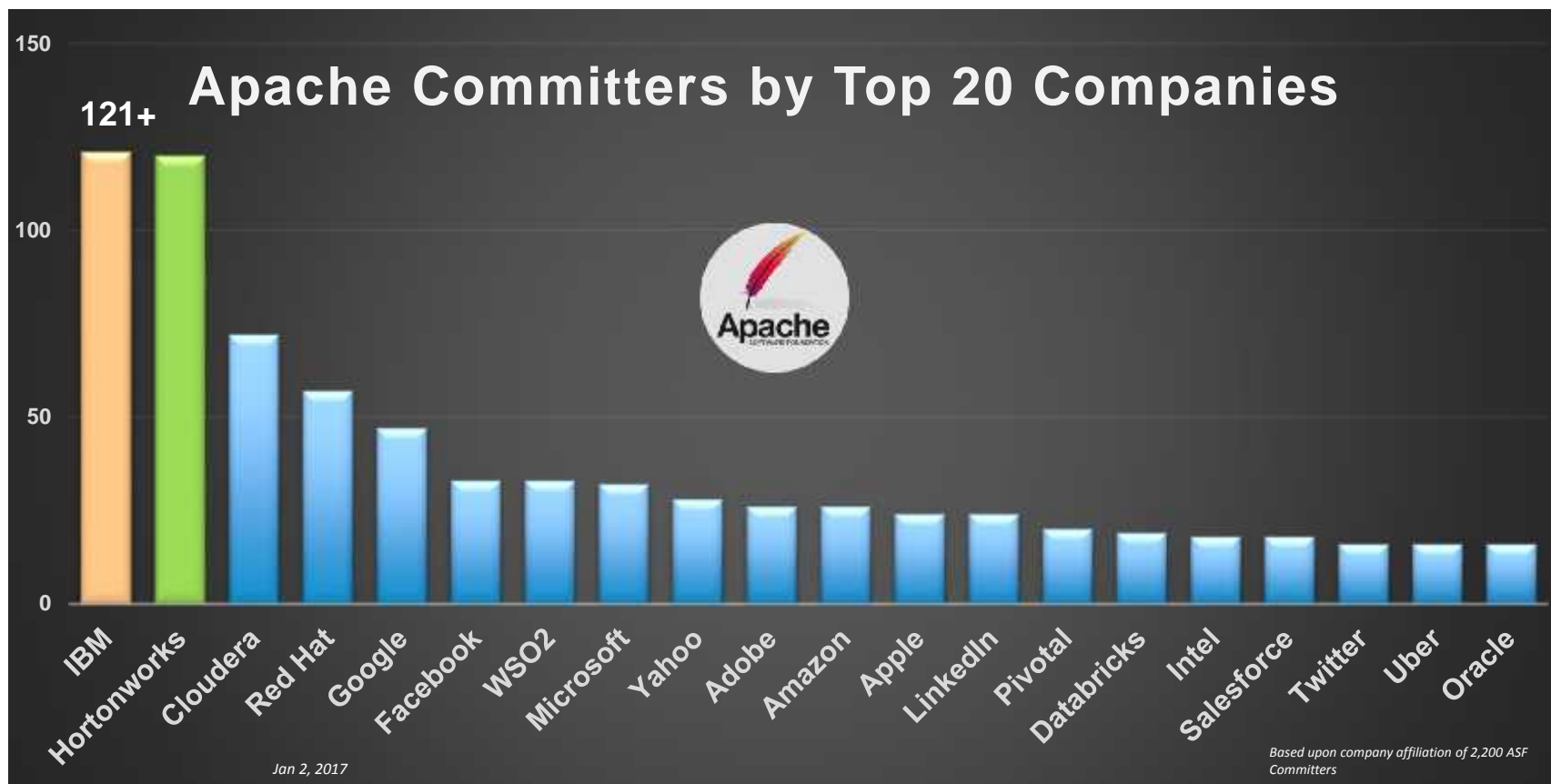
**IBM** + **Hortonworks**

- **#1 Rank by Gartner** 2017 Data Science Magic Quadrant
- **Leader in SQL technology** for Hadoop (www.tpc.org)
- **Leader in data and analytics** solutions for hybrid cloud
- Provides Data Science **& Machine Learning**

- **Leader in Hadoop** Open Source Distribution
- **1000+ customers** and 2100+ ecosystem partners
- **Hadoop original architects, developers** employed by Hortonworks
- **Provides Open Hadoop** Data Platform

**Commitment to progressing advanced analytics through open source**

# IBM and Hortonworks - Open Source Commitment

**…and our combined commitment to Open Standards is Unmatched.**

# IBM Big Data High Value with Hortonworks



IBM's Offerings Unlock the value of Hadoop Data

**IBM BigIntegrate / BigQuality / BigMatch**
- Large scale data ingest & transformation
- Data analysis, cleansing, & monitoring
- Accurate linkage of customer data

**#1 Data Science Platform: DSX**
- Community and social features to provide collaboration
- The best of open source and IBM value-add to create state-of-the-art data products
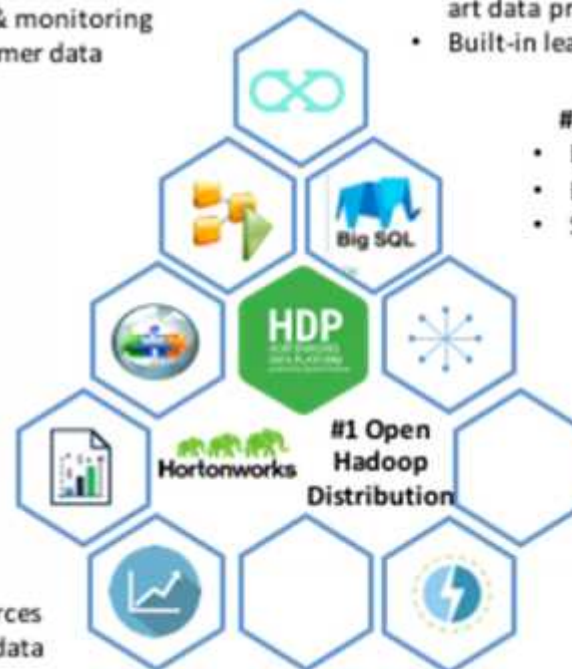- Built-in learning and advanced tutorials

**IBM Information Governace Catalog**
- Understand, Curate, and Govern
- Business level glossary and Catalog
- Comprehensive data lineage and tool impact analysis

**#1 SQL Engine for Hadoop: Big SQL**
- Data virtualization layer
- Large data volume, extremely complex query support
- Supports low latency, high concurrency workloads

**Cognos, Watson Analytics**
- Self service analytics capabilities
- Guided Analytics Discovery
- Natural Language Dialogue

**IBM Big Replicate / IBM Data Replication**
- Multiple Hadoop distributions to Hadoop
- Source Application to Hadoop Replication
- Provides HA/DR, with virtually zero RTO/RPO
- On-Prem to Cloud and Cloud to On-Prem

**SPSS**
- Further embrace and extend Open source
- Integrate with other IBM offerings & data sources
- Energize your Analytics (text analytics for Big data on System-T)

**IBM Streams**
- Built-in streaming analytics
- Open architecture. Built for Speed
- Integrated Dev Environment

HDP

Hortonworks

#1 Open Hadoop Distribution

Big SQL

# Db2 Big SQL – For all WH Needs in Hadoop

**SQL-based Application**

**Common SQL Engine Client driver**

**Big SQL Engine**

**SQL MPP Run-time**

**Data Storage**

**DFS**

**Hadoop**

**SQL**

Ad-hoc queries, data preparation

Federation

Operational with fast lookups

High performance and scalability

Integrated Spark and Machine Learning

Complex SQL, Deep analytics, Many users

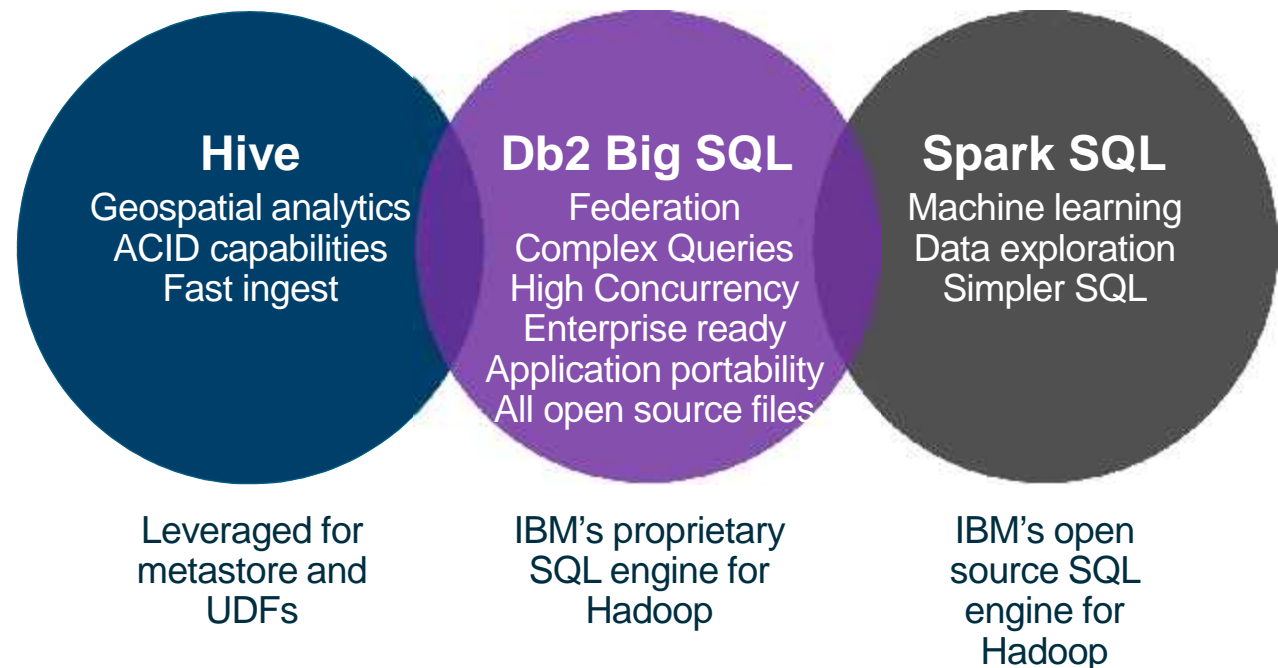Application portability

# Db2 Big SQL V5.0

| Applications | •ETL<br>•Reporting<br>•Data mining<br>•Deep analytics | •Reporting<br>•Complex queries<br>•BI Tools: Cognos, Tableau, etc | •Query EDW<br>•Join data<br>•Use ML | •Reuse applications<br>•Reuse skills | •Ad-hoc, exploratory<br>•BI tools: Cognos, Tableau, etc |
|---|---|---|---|---|---|
| Capabilities | Batch SQL<br>(minutes to hours) | Interactive SQL<br>(seconds to minutes) | Data augmentation<br>(Spark integration) | Application portability | Self-service / Interactive BI<br>(Sub-second) |
| Core | SQL compatibility – Db2, Oracle, Netezza | SQL and NoSQL Structured & Unstructured | DSM, Ambari | MQTs | Ranger |
| | Advanced cost-based optimizer | Federation | Automatic memory management | Elastic boost – logical worker nodes | Roles |
| | Comprehensive ANSI SQL coverage | Spark Integration | Automatic workload management WLM | Query rewrite for optimized execution | SQL based RBAC |
| | Core SQL Engine | Integration | Administration | Performance | Security |

www.tpc.org – check out TPC-H and TPC-DS – Big SQL vs Impala vs Hive
Db2 Big SQL 5.0 is **2X** faster than Hive LLAP with Tez – and much more functional
Db2 Big SQL 5.0 is **3X** faster than Spark SQL 2.1

# Combining Hadoop Technologies

Not Mutually Exclusive. Hive, Db2 Big SQL & Spark SQL can co-exist and **complement** and **leverage** each other in a cluster

**Hive**
Geospatial analytics
ACID capabilities
Fast ingest

**Db2 Big SQL**
Federation
Complex Queries
High Concurrency
Enterprise ready
Application portability
All open source files

**Spark SQL**
Machine learning
Data exploration
Simpler SQL

Leveraged for metastore and UDFs

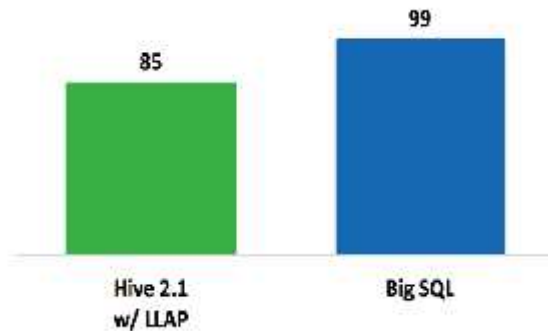IBM's proprietary SQL engine for Hadoop

IBM's open source SQL engine for Hadoop

# Query Performance at a Glance – vs Hive LLAP with Tez

## HADOOP-DS @ 10TB
**85 COMMON QUERIES**

**WORKING COMPLIANT QUERIES: 6-streams**



**PERFORMANCE: 6-streams**

## Db2 Big SQL **2.3x** FASTER



**WORKLOAD**
SCALE FACTOR: **10 TB**
FILE FORMAT: **ORC** (ZLIB)
CONCURRENCY: **6 STREAMS**
QUERY SUBSET: **85 QUERIES**

**INTERESTING FACTS**

**FASTEST QUERY**
**5.4x** FASTER (Db2 Big SQL: 1.5 SEC, HIVE: 8.1 SEC)

**SLOWEST QUERY (QUERY 67)**
**1.7x** FASTER (Db2 Big SQL: 6827 SEC, HIVE: 11830 SEC)

**Db2 Big SQL** FASTER FOR **80%** OF QUERIES RUN

**PERFORMANCE: 1-stream**

## Db2 Big SQL **1.8x** FASTER



**STACK**
HDP 2.6.1
Db2 Big SQL 5.0.1
HIVE 2.1 LLAP ON TEZ

## RESOURCE UTILIZATION:
**6-STREAMS**

**1.5x** FEWER CPU CYCLES USED

# Query Performance at a Glance – Db2 Big SQL & Spark SQL

*Leads performance metrics on high volumes of data and concurrent streams*
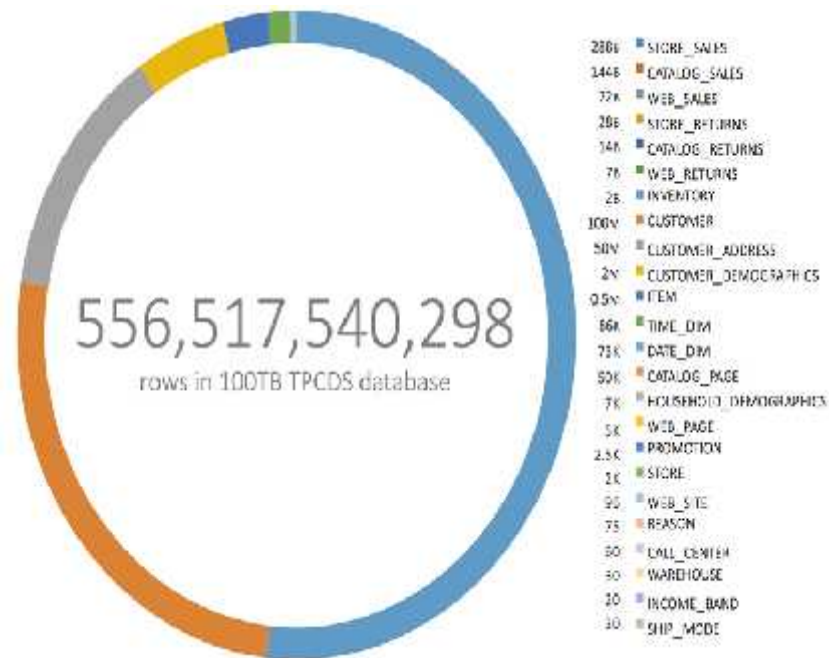
## SNAPSHOT OF 100TB HADOOP-DS



556,517,540,298
rows in 100TB TPCDS database

| | |
|---|---|
| 280B | STORE_SALES |
| 144B | CATALOG_SALES |
| 72B | WEB_SALES |
| 28B | STORE_RETURNS |
| 14B | CATALOG_RETURNS |
| 7B | WEB_RETURNS |
| 2B | INVENTORY |
| 100M | CUSTOMER |
| 50M | CUSTOMER_ADDRESS |
| 2M | CUSTOMER_DEMOGRAPHICS |
| 0.5M | ITEM |
| 86K | TIME_DIM |
| 73K | DATE_DIM |
| 50K | CATALOG_PAGE |
| 7K | HOUSEHOLD_DEMOGRAPHICS |
| 5K | WEB_PAGE |
| 2.5K | PROMOTION |
| 1K | STORE |
| 90 | WEB_SITE |
| 75 | REASON |
| 60 | CALL_CENTER |
| 30 | WAREHOUSE |
| 30 | INCOME_BAND |
| 30 | SHIP_MODE |

### PERFORMANCE
Db2 Big SQL 5.0 is **3.2x** faster than Spark SQL 2.1
*(4 Concurrent Streams)*

Big SQL — **13.7 hours**

Spark SQL — **43.2 hours**

### COMPRESSION
**60%**
**SPACE SAVED WITH PARQUET**

### AVERAGE CPU USAGE
**76.4%**

### MAX I/O THROUGHPUT
**READ** 4.4 GB/SEC
**WRITE** 2.8 GB/SEC

### WORKING QUERIES
83 — Spark SQL
99 — Big SQL

### I/O (vs Spark)
Db2 Big SQL reads **12x** less data
Db2 Big SQL writes **30x** less data

Blog on benchmark: https://developer.ibm.com/hadoop/2017/02/07/experiences-comparing-big-sql-and-spark-sql-at-100tb/

# Federation – Query one connection and Virtualize Heterogeneous Data

**Db2 Big SQL queries heterogeneous systems in a single query**

**Only SQL-on-Hadoop that virtualizes more than 10 different data sources: RDBMS, NoSQL, HDFS or Object Store**

**Transparent**
- **Appears to be one source**
- **Programmers don't need to know how / where data is stored**

**High Function**
- **Full query support against all data**
- **Capabilities of sources as well**

**Autonomous**
- **Non-disruptive to data sources, existing applications, systems.**

**High Performance**
- **Optimization of distributed queries**

| MS SQL Server | Netezza (PDA) | Oracle | PostgreSQL | Teradata | DB2 LUW, Db2z, DB2 on i | Informix |

Db2 Big SQL

Spark

NoSQL

| WebHDFS | Object Store (S3) | Hive | HBase | HDFS |

Hortonworks Data Platform (HDP)

ML Model

# Federation: Rich Capabilities that Brings Data Together

✓ *Easily access information on demand*

✓ *Combine data in Hadoop with disparate sources to form a data lake*

✓ *Quickly extend your data warehouse by enriching it*



| Connect | Query | Monitor | Data Placement |
|---|---|---|---|
| ▪ Quick access to Data value<br>▪ Common Framework<br>▪ ODBC/JDBC<br>▪ Spark integration enables new data sources<br>▪ Connect all data sources in single query | ▪ Intelligent Query Routing<br>▪ Cost-based optimizer<br>▪ SQL pushdown<br>▪ Local data caching<br>▪ ANSI-compliant SQL | ▪ Easily define & manage through a common UI<br>▪ Simple point & click to discover and query<br>▪ Monitor and visualize active queries | ▪ Schema conversion when moving data<br>▪ Bulk data copy to Hadoop<br>▪ Filtered subsets of data |

# Application Portability: Move Applications without Re-tooling



Data warehouse offload to Hadoop is now made easy:

- Write one, run anywhere…
- Easy porting of applications
- Reuse skills of DBAs/ developers who know ANSI SQL

Db2 Big SQL is the best platform for offloading Oracle Data Marts and Warehouses to Hadoop

# Oracle Compatibility - SET sql_compat='ORA'

| DB2, Big SQL | Oracle, Netezza, Postgres |
|---|---|
| TRANSLATE (expr, to_str, from_str) | TRANSLATE (expr, from_str, to_str) |

## Same function, parameters reversed!

- **SQL_COMPAT global variable enables support for both parameter orderings (and other syntax/behavioral conflicts when offloading SQL to Hadoop)**

- **Excellent Oracle PL/SQL support! (New for V5.0)**
  - SQL data-access-level enforcement
  - Enforce data access levels at run time rather than at compile time.
  - Oracle database link syntax (@ symbol)

- **Note:**
  - Setting of the DB2_COMPATIBILITY_VECTOR registry variable (inherited from DB2) is **not recommended** in Big SQL. Custom compatibility features should be enabled **only** by using the SQL_COMPAT global variable.

# Oracle PL/SQL Support

```
set sql_compat='ORA'        ← Easy session variable to switch modes!

create or replace procedure plsql_proc (fetchval out integer)
as

cursor cur1 is
select count(*) from syscat.tables ;

-- begin

open cur1 ;

fetch cur1 into fetchval ;
close cur1 ;

end
```

Big SQL is the best platform for offloading Oracle Data Marts and Warehouses to Hadoop

- Big SQL V4.2 already supports some Oracle SQL compatibility (but not PL/SQL)
- Big SQL V5.0 adds support for Oracle PL/SQL procedural language

# Query Execution

*Here's why Db2 Big SQL can get you the best execution for complex queries and many concurrent users with high performance*
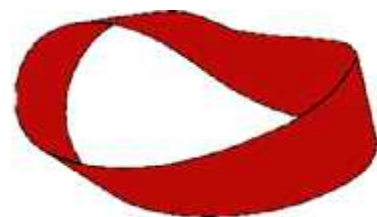
## *Performance*

**Materialized Query Tables**

**Advanced Statistics**

**Elastic Boost**

**World Class Cost Based Optimizer**

## *Concurrent users*

**Self Tuning Memory Manager**

**Advanced Workload manager**

## *Complex query*

**SQL Compatibility**

**Hardened runtime**

**Query rewrite**

**Native Row & Columnar stores**

# Big SQL – Rich Analytics

**Big SQL** is a powerful hybrid analytical engine

Offers leading performance metrics on high volumes of data to combine, transform, cleanse data in a secure environment, to generate a dataset to derive insights on data

Data engineer



Unstructured data

Structured data

Semi-structured data

**Big SQL**

data integration

data security

Filter    Aggregate

Union    Join

Expression

What-If Analysis

Jupyter Notebooks

Cognos Analytics

REPORTS

**Data input**          **Generate quality data**          **Business optimization**

# Self-service Analytics: Democratize Data Science and ML

*Leverage **Db2 Big SQL** throughout your journey*

| Data Ingestion | → | Data Transformation/ Data Science/ Machine Learning | → | Data Visualization |
|---|---|---|---|---|
| *Virtualize disparate data sources like Hadoop, RDBMS, and Object Stores (S3) to join data in a single query* | | *Manipulate data and operationalize data science models written in various languages* | | *Perform data discovery, analyze, and visualize business results in notebooks or other BI tools* |

# Operationalize Machine Learning Models using SQL



For more details check the blog: https://developer.ibm.com/hadoop/2017/11/07/ibm-big-sql-machine-learning-demo/

# Big SQL - Security

Big SQL is the **most secure analytical engine** that offers row and column level access control (RCAC) among other security settings

Data engineer

**Role Based Access Control enables separation of Duties / Audit**

BRANCH_A

BRANCH_B

FINANCE (security admin)

## Row Level Security

| EMPNO | FIRST_NAME | SALARY | BRANCH_NAME |
|---|---|---|---|
| 2 | Chris | 29007.57 | Branch_A |
| 3 | Paula | 14987.06 | Branch_A |
| 5 | Pete | 19114.22 | Branch_A |
| 8 | Chrissie | 24922.36 | Branch_A |

*Total: 8 Selected: 0*      < 1 >      10 | 25 | 50 | 100

## Row and Colum Level Security

| EMPNO | FIRST_NAME | SALARY | BRANCH_NAME |
|---|---|---|---|
| 1 | Steve |  | Branch_B |
| 4 | Craig |  | Branch_B |
| 6 | Stephanie |  | Branch_B |
| 7 | Julie |  | Branch_B |

| EMPNO | FIRST_NAME | SALARY | BRANCH_NAME |
|---|---|---|---|
| 1 | Steve | 25970.38 | Branch_B |
| 2 | Chris | 29007.57 | Branch_A |
| 3 | Paula | 14987.06 | Branch_A |
| 4 | Craig | 22518.93 | Branch_B |
| 5 | Pete | 19114.22 | Branch_A |
| 6 | Stephanie | 26183.81 | Branch_B |
| 7 | Julie | 13829.91 | Branch_B |
| 8 | Chrissie | 24922.36 | Branch_A |

*Total: 8 Selected: 0*      < 1 >      10 | 25 | 50 | 100

# Big SQL and Apache Ranger Integration



- **Setup ACLs for access to Big SQL tables:**
  - create, alter, analyze, load, truncate, drop, insert, select, update, and delete.
- **Supports Ranger Audit**
  - Big SQL access audit records written to HDFS and/or Solr
- **If also using Ranger Plugin for Hive – operates independent of Big SQL plugin**

# Big SQL Tables over S3 Object Storage

> CREATE HADOOP TABLE staff ( … )
> LOCATION
>        '**s3a://**s3atables/staff';

- **Create Tables over Data residing in Object Store directly (no copy required into Hadoop)**

- **Once configured, Object Store tables work like any other table in Big SQL**

- **Benefits:**
  - No need to copy data into Hadoop first! Query data where it resides.
  - Partitioning supported!

- **Tradeoff:**
  - Expect reduced performance relative to HDFS local tables

LOAD FROM
Object Store
also supported!

IBM Corporation

# Big SQL Tables over WebHDFS (Technical Preview)

CREATE HADOOP TABLE staff ( … )
    PARTITIONED BY (JOB VARCHAR(5))
LOCATION 'webhdfs://namenode.acme.com:50070/path/to/table/staff';

WebHDFS
REST call

Big SQL

*Local Hadoop Cluster*

*Remote Hadoop Cluster
or
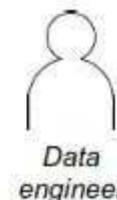WebHDFS enabled Storage*

- **Transparently access data on any platform implementing WebHDFS**
    - Examples: Microsoft Azure Data Lake (ADL) service

- **Once setup, WebHDFS tables work like any other table in Big SQL**

- **Technical Preview Limitations:**
    - WebHDFS via Knox not supported
    - Performance not well understood. Reduce performance expected.

LOAD FROM
WebHDFS
also supported!

# Big SQL – Integration with Yarn and Spark

**Big SQL is a self-tuning memory management SQL engine that integrates with Spark 2.1**

Data engineer

Share data in memory

Big SQL worker

Spark executor

Spark 2.1 is a powerful analytic co-processor that complements the rich SQL functionality of Big SQL

Bi-directional integration allows Spark jobs can be executed from Big SQL

Tight integration with Spark enables Big SQL worker and Spark Executor to communicate in memory without writing to disk

Big SQL

Launch

HDFS

APACHE Spark™

# Big SQL – The ONLY engine with Deep Integration with Spark

# Exploit Big SQL from Spark

```
import org.apache.spark.sql.Dataset;
import org.apache.spark.sql.Row;
…
Dataset<Row> tableDf = sqlCtx.read()
  .format("jdbc")
  .option("driver", "com.ibm.db2.jcc.DB2Driver")
  .option("url", "jdbc:db2://server1.foo.bar.com:32051/BIGSQL")
  .option("user", "joe")
  .option("password", "joespwd")
  .option("dbtable", "myshcema.mytable")
  .load();


tableDf.createOrReplaceTempView("myTable");
Dataset<Row> queryDF =
    spark.sql("SELECT col2, col3 FROM myTable WHERE col1 > 100");
```

Big SQL secures data for self-service data exploration.

Used this way, Spark users are subject to Big SQL row/column security

- ▪ **Requirements:**
  - – db2jcc.jar must be added to the classpath of the Spark application (found in /home/bigsql/java/)

# Exploit Spark from Big SQL
## *Example: Spark Schema Discovery for JSON*

```
SELECT doc.*
FROM TABLE(
        SYSHADOOP.EXECSPARK( class => 'DataSource',
                                load => 'hdfs://host.port.com:8020/user/bigsql/demo.json')
) AS doc
WHERE doc.language = 'English';
```

Structure of JSON document determined at run time

- **Bring the best of Spark into Big SQL!**

  – Machine Learning

  – Cache remote tables (Spark has rich library of connectors)

  – Graph Processing

  – General in memory processing

# Apache Slider

- **Apache Slider**
  - Enables long running services (e.g. Big SQL) to integrate with YARN (similar to HBase)
  - Provides:
    - Implementation of Application Master
    - Monitoring of deployed applications
    - Component failure detection and restart capabilities
    - Flex API for adding/removing instances of components of already running

- **Apache Slider does not yet have a GUI nor Ambari integration.**
  - Big SQL operations for Slider can be executed through two methods:
    - Big SQL Service Actions in Ambari
    - Command line scripts

# Big SQL + YARN Integration
**Dynamic Allocation / Release of Resources**

Users

Big SQL
Head

Slider Client

YARN
Resource
Manager
& Scheduler

*Big SQL Slider package
implements Slider Client APIs*

YARN
components

Slider
Components

Big SQL
Components

Container

NM | NM | NM | NM | NM | NM

Big SQL
AM

Big SQL
Worker

*Stopped workers
release memory
to YARN for
other jobs*

Big SQL
Worker

*Stopped workers
release memory
to YARN for other
jobs*

*Stopped workers
release memory
to YARN for
other jobs*

Big SQL
Worker

HDFS

74

# Big SQL Elastic Boost –
# Multiple Workers per Host
*More Granular Elasticity*

Up to 50% more query performance

Users

Big SQL Head

Slider Client

YARN Resource Manager & Scheduler

YARN components

Slider Components

Big SQL Components

Container

NM

NM

NM

NM

Big SQL AM

NM

NM

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

Worker

HDFS

# Elastic Boost Improves INSERT Performance

INSERT…SELECT performance with Elastic Boost



In each scenario, the same TOTAL CPU/memory is used

- **For both 1 and 10 TB TPC-DS dataset**
  - 2 Workers/Node: **1.6x speedup**
  - 4 Workers/Node: **2.2x speedup**

# Big SQL 5.0 – How it fits with Hortonworks

- **Big SQL deploys on top of Hortonworks Data Platform(HDP)**
  - *Includes:* IBM Support for Big SQL
- **Hortonworks Data Platform for IBM (Support only)**

- **Hortonworks Data Platform can be downloaded for FREE.**

| IBM Big SQL | Hortonworks Data Platform for IBM |
|---|---|
| includes support for Big SQL | Support Offering for HDP |
| BUY | BUY |

**Hortonworks Data Platform**

FREE

Spark

**Capability**
Apache Hadoop and
Apache Spark and Ecosystem

**Support**
Community Support

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**

- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# Event-Driven Systems Span Many Industries



Multi-channel customer sentiment and experience a analysis

Detect life-threatening conditions at hospitals in time to intervene

Predict weather patterns to plan optimal wind turbine usage, and optimize capital expenditure on asset placement

Make risk decisions based on real-time transactional data

Identify criminals and threats from disparate video, audio, and data feeds

# Industry Use Cases

### Retailer Loyalty Program

Integrate streamed payment, couponing events, climate, calendar, mobile data. measure refine, deliver better couponing and loyalty system

### Smart Metering/Smart Grid

Deliver a Integrated platform for optimizing energy usage, capacity and billing across a smart grid system

### Banking Risk Exposure

Combine account transactions from across the bank to provide a master ledger for real-time risk exposure and fraud identification

### Satellite Tracking System

Track satellites in real time  and produce analytics on operations and performance

### Intelligent Manufacturing

Deliver real-time monitoring framework for automated production lines, providing productivity, preventive maintenance, and reporting

### Transactional to Analytics Consolidation

Capture your transactions and augment with external data into an analytics platform for deeper analytics

# What is Db2 Event Store?

*A unified offering for Fast Data which delivers…*

IBM **Db2 Event Store**

**.1** **Lightning Fast Ingest**

- 1 Million inserts per second per node
- Ingest scales linearly with added nodes
- Data ingested quickly, then refined and enriched

**2** **Real-time Analytics**

- Real-time analytics over ALL ingested data
- Super-fast lookups and intelligent scans
- Integrated machine learning capabilities

**3** **Integrated and Highly Available**

- Packaged and integrated with IBM Data Science experience; available Streams sink
- Remains available on node failure
- Architected to scale to very large clusters

**4** **Built for Data Sharing and Efficiency**

- Writes to shared storage in Parquet format
- Able to leverage low-cost object storage
- Single copy of the data

**Parquet**

# Db2 Event Store
## *Integrated System for Managing Events*

IBM Event Store

## Simple to Deploy and Scale Container Delivery

IBM Data Science Experience

Native Access to
Spark Streaming    kafka    Scala    Java

**IBM Project Event Store**    Spark

IBM Streams

Columnar Organized

In-Memory Indexing

Replication For High Availability

Fast Performance

IBM BigSQL

Parquet    Open Apache Parquet Data Format

# Db2 Event Store – Competitive Positioning

**Db2 Event Store provides everything in the box**

**Reduced architectural components
Docker Container Delivery
Open Data Access**

**Complex Manual Architecture**

**Put together your own open source components
Not everything works together
Hard to maintain**



**Simplified Approach With IBM Db2 Event Store**

*VS.*

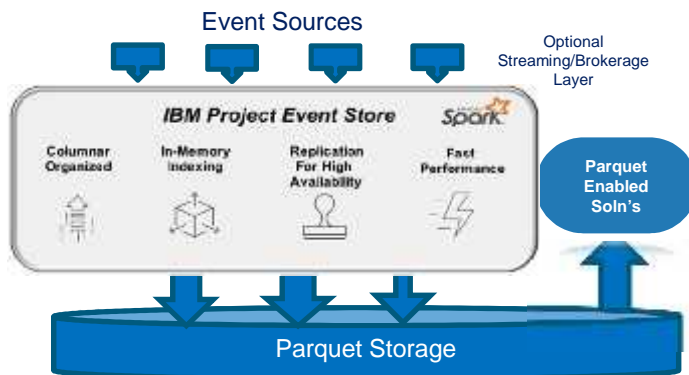**Digital Company Example**

# Db2 Event Store

# Db2 Event Store:  Unified Data Access and Management

**Spark / Event Store API**

**Data Server Manager**
*Unified Management platform*

**Data Server Manager Client**

**IBM Project Event Store** Spark

| Columnar Organized | In-Memory Indexing | Replication For High Availability | Fast Performance |

IBM DB2

Native DRDA Federated Access

BigSQL Platform

Native parquet MPP Access

Unifying Event Data with traditional Stores

Object Storage / HDFS Parquet Data format

# Db2 Event Store: Architecture

## *Understanding the Engine and Components*

IBM Event Store

**High Speed Ingest**
kafka  IBM Streams  Scala  Spark Streaming

**Real-Time Insights**  Machine Learning
Scala  Java  Spark

**IBM BIGSQL  Parquet Compatible tools**
Spark

**IBM Event Store Cluster**

Event Data Management Engine
Columnar Organized  In-Memory Indexing  Replication For High Availability  Fast Performance

Event Data Management Engine
Columnar Organized  In-Memory Indexing  Replication For High Availability  Fast Performance

Event Data Management Engine
Columnar Organized  In-Memory Indexing  Replication For High Availability  Fast Performance

Highly Available Distributed Storage    Parquet    GlusterFS    cleversafe    Open Data Format In-Memory Data grid

# Db2 Event Store

**Demo**

# Topics for Today

- **Strategy Overview**

- **Db2 V11.1.3.3 – Introduction !!**

- **Private Cloud – Introduction !!**

- **Flex Points and HDM Offering**

- **Appliance News**
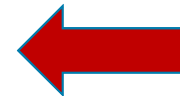
- **Hadoop and Open Source**

- **Event Processing**

- **Next Generation Data Virtualization**

# Analytics Today…



Data Lake or Data Warehouse

- Costly and Complex

- High Latency to copy and synchronize

- Available compute resources under-utilized

- Error prone and difficult to retain data integrity

# IBM Queryplex
## *An emerging technology now in beta trial*

**Queryplex Constellation**

### Query anything, anywhere.

**1** Query **many diverse data sources** across cloud, on-premise and mobile with advanced analytics using the most popular languages and tool

SQL, Spark, R, Notebooks, Python, Data Science Experience (DSX), Cognos Analytics, common Analytics tools
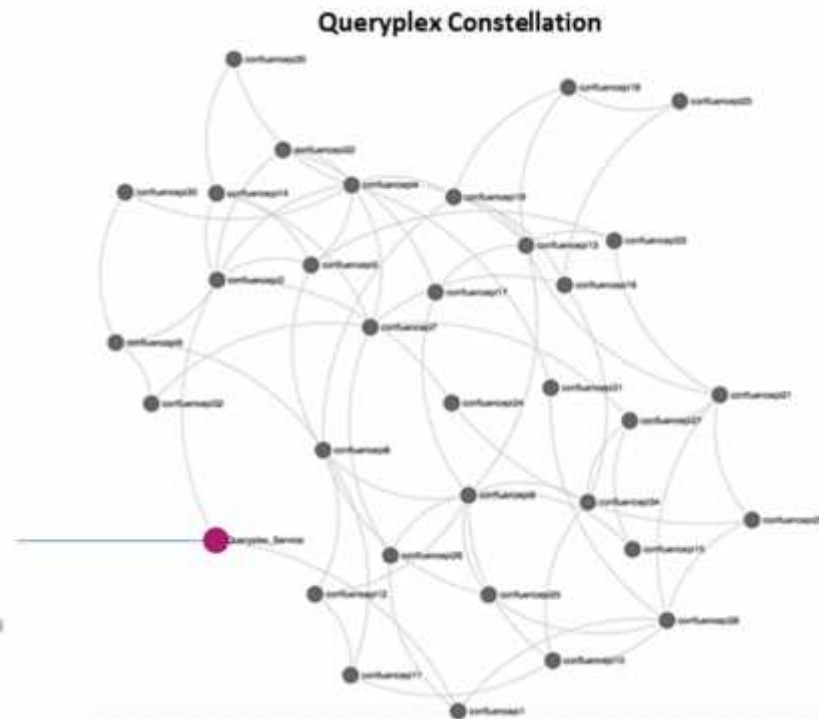
Analytics Application

### Query many sources as one with extreme simplicity.

**2** Connect **few to many devices and data stores** into a single self balancing constellation. Avoid the complexity of centralized copies. Data only persists at the source.
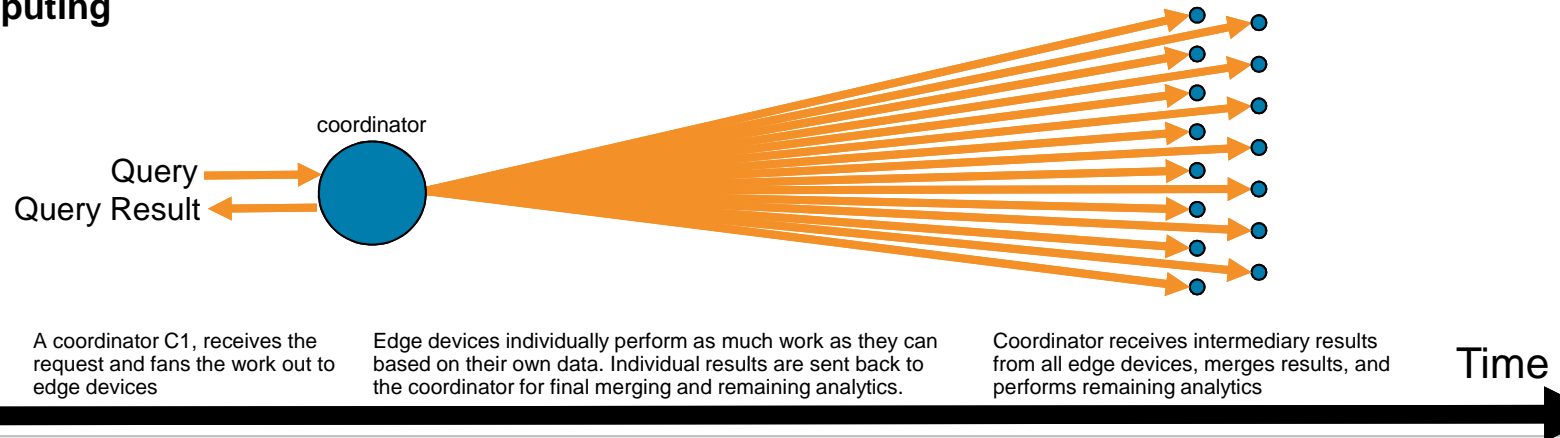
### Massive speedup.
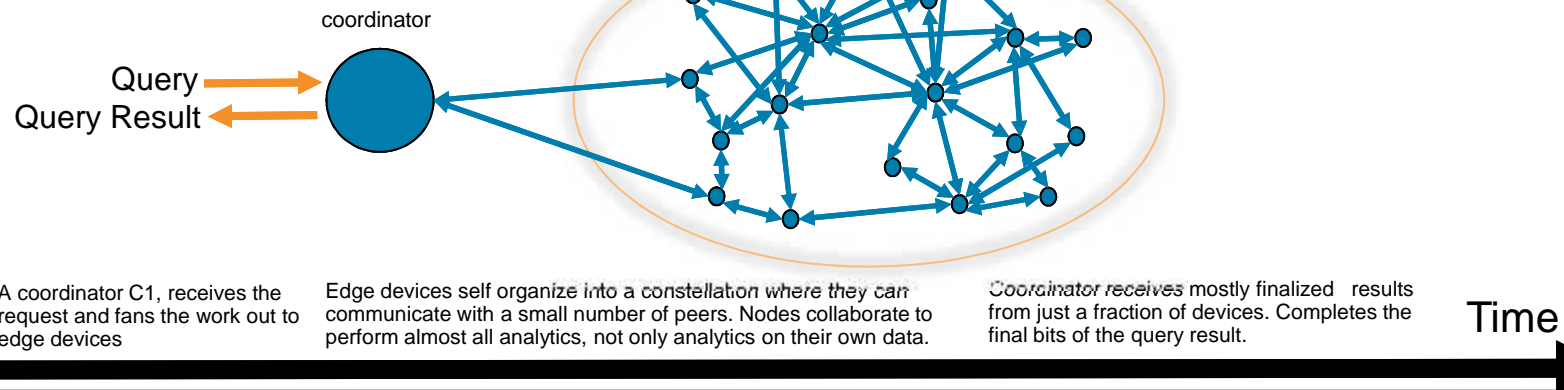
**3** **Many times acceleration** using the power of every device.

# IBM Queryplex's Computational Mesh

## Edge Computing



coordinator

Query

Query Result

| Query issued against the system | A coordinator C1, receives the request and fans the work out to edge devices | Edge devices individually perform as much work as they can based on their own data. Individual results are sent back to the coordinator for final merging and remaining analytics. | Coordinator receives intermediary results from all edge devices, merges results, and performs remaining analytics |

Time

## Queryplex's Computational Mesh



coordinator

Query

Query Result

| Query issued against the system | A coordinator C1, receives the request and fans the work out to edge devices | Edge devices self organize into a constellation where they can communicate with a small number of peers. Nodes collaborate to perform almost all analytics, not only analytics on their own data. | Coordinator receives mostly finalized results from just a fraction of devices. Completes the final bits of the query result. |

Time

# IBM Queryplex - Supported Languages & Data Sources

| Query Languages | |
|---|---|
| SQL (ANSI) | ✓ |
| SQL (Oracle) | ✓ |
| SQL (DB2) | ✓ |
| SQL (PostgreSQL, Netezza) | ✓ |
| Scala | ✓ |
| PL/SQL | *Future* |
| SQL PL | *Future* |
| PySpark | ✓ |
| Python | ✓ |
| R & SparkR | ✓ |

| Mix Any Combination of Data Sources | | | |
|---|---|---|---|
| Oracle | ✓ | Excel | ✓ |
| DB2 | ✓ | CSV (delimited text) | ✓ |
| Netezza | ✓ | MongoDB | ✓ |
| PostgreSQL | ✓ | Accumulo | *Future* |
| Informix | ✓ | Redis | *Future* |
| MySQL | ✓ | Cloudant | *Future* |
| SQLServer | ✓ | | |
| DerbyDB | ✓ | | |

# IBM Queryplex - Potential Use Cases

| Industry | Use Case |
|---|---|
| Telco | 5G Wireless and Enterprise IoT (Devices anywhere) |
| Telco | Cell tower and site monitoring for Operations and Maintenance |
| Telco | Cell site subscriber metadata analytics for Law Enforcement |
| Telco | Set Top Box home applications, monitoring, Content access statistics |
| Energy & Utilities | Distribution network monitoring and maintenance |
| Energy & Utilities | Smart metering |
| Manufacturing & Cross/Enterprise | Time sensitive data queries |
| Insurance | Auto usage device monitoring |
| Cross/Enterprise | Data Virtualization |
| Cross/Enterprise | Data provisioning to untrusted external entities |
| Gaming | Real-time gaming queries |
| Media & Entertainment | Subscriber viewing and content correlation |
| Military | IoT Sensors |

# IBM Queryplex – Interested in hearing more ?



IBM Queryplex
The power of many together

http://queryplex.com

**Les King**

**Director, Hybrid Data Management Solutions**

**May, 2018**

lking@ca.ibm.com

ca.linkedin.com/pub/les-king/10/a68/426

# Hybrid Data Management Strategy and New News !