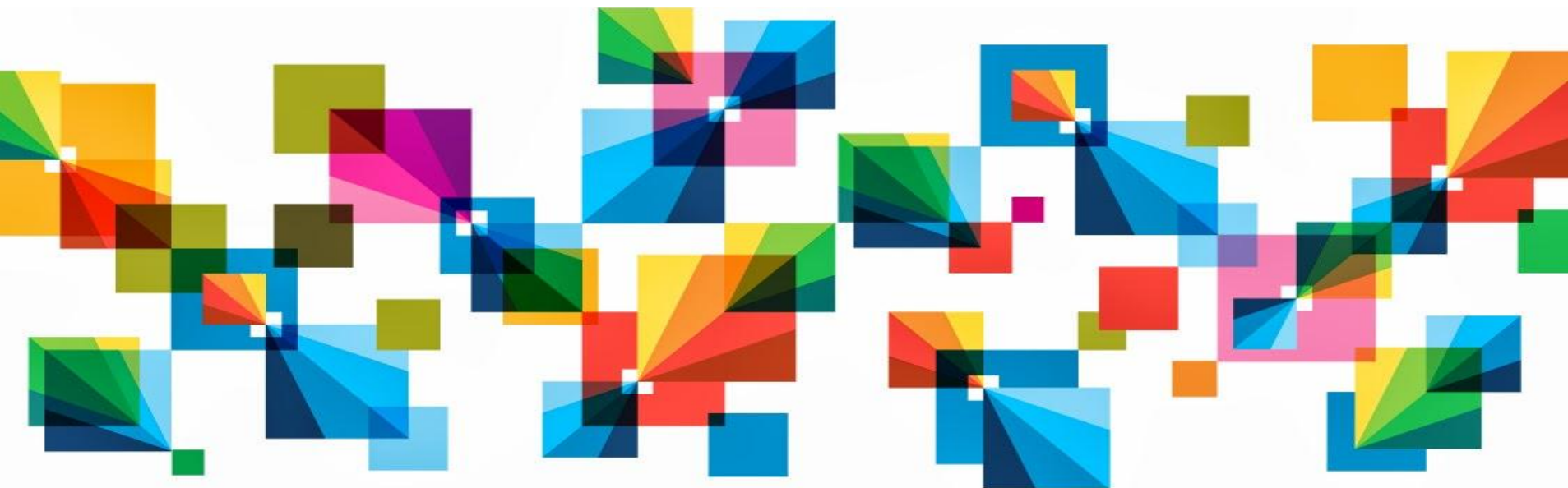


Amplia tu Estrategia de Análisis

Modernización del Almacén de Datos



SOF Quiere Ahorrar Dinero y Cumplir con sus Requisitos Regulatorios

Tengo una gran cantidad de datos históricos en mi almacén de datos relacional que tengo que mantener por razones regulatorias.

Me gustaría que hubiera un lugar más económico para almacenar esos datos.



CTO

Debe considerar la construcción de un almacén de datos moderno. Le voy a enseñar cómo.

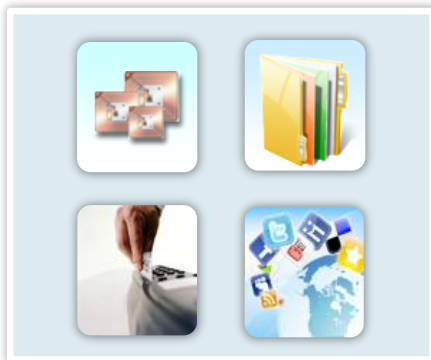


IBM

Modernización de Almacenamiento de Datos Definido



Integrar las tecnologías de “big data” con el entorno de su almacén de datos para obtener nuevas perspectivas al negocio y optimizar su infraestructura de almacén de datos.



Nuevos tipos de información

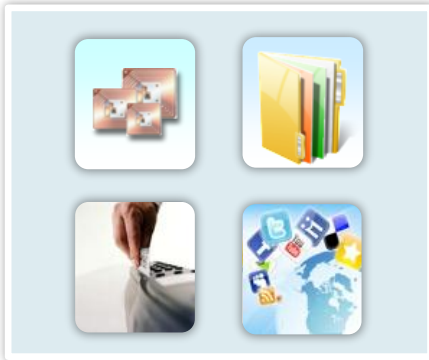


Optimice la infraestructura del almacén de datos existente



Reduzca el costo total de propiedad del almacén de datos

La Modernización del Almacén de Datos Aborda Muchos Desafíos



Nuevos tipos de información

- Aproveche nuevos fuentes de datos para perspectivas no antes posible
 - Internet de cosas
 - Datos sociales
- Integración con datos en movimiento
- Perspectivas en minutos y horas no días y semanas



Optimización de la infraestructura AD existente

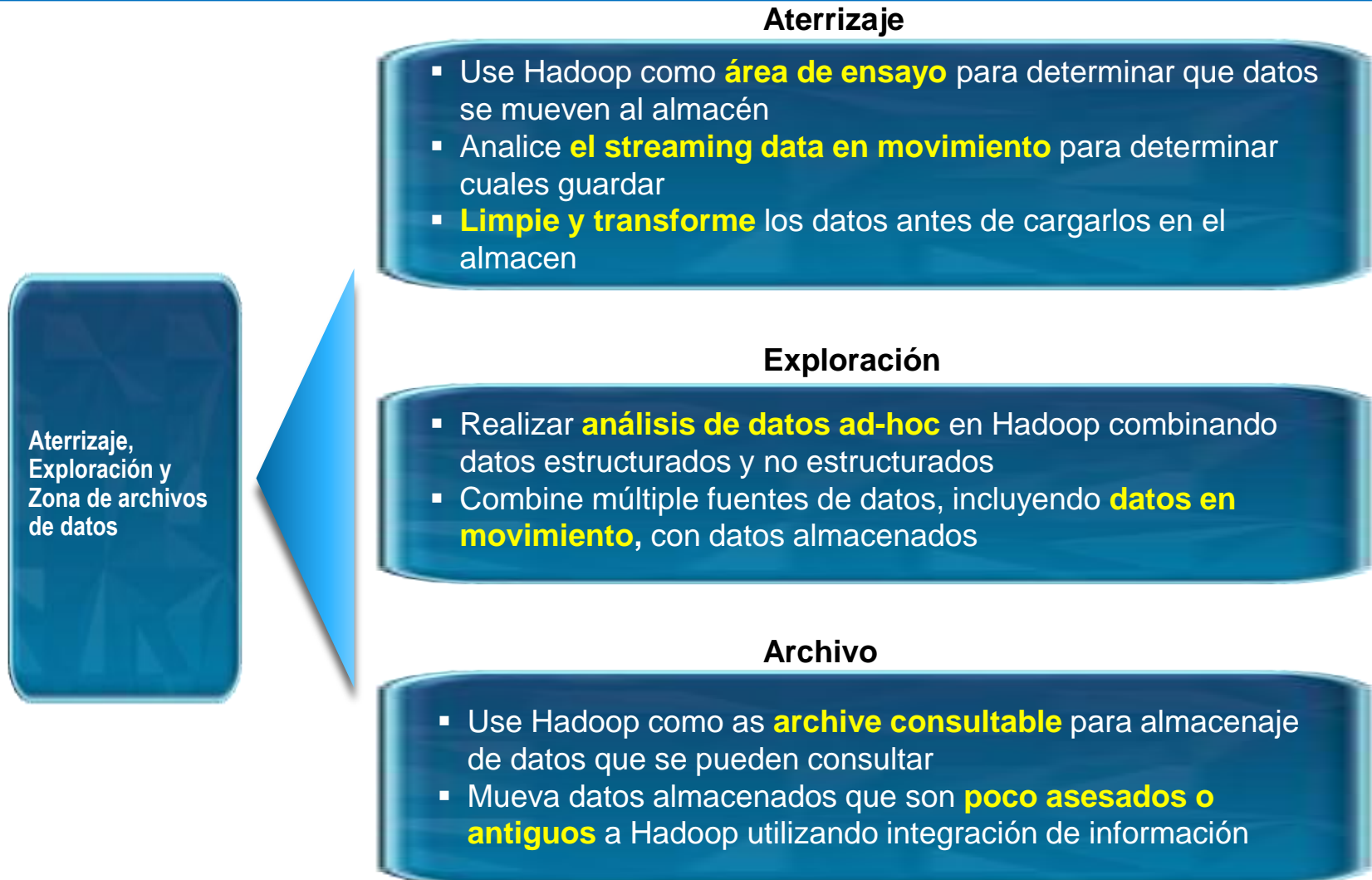
- Optimice la infraestructura del almacén mediante la descarga del pre-procesamiento
- Optimice el almacenamiento de datos a través del almacén y Hadoop
- Mas rápido ciclo de vida de informes
- Visión empresarial comprensivo de todos los datos



Reduce el costo total de propiedad del AD

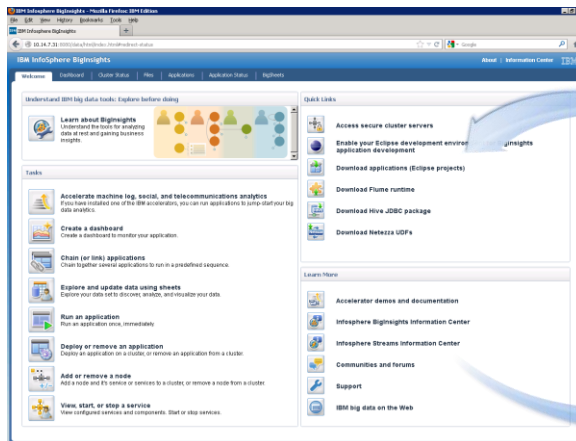
- Migrar datos rara vez usados a Hadoop
- Incremente la optimización del HW del AD alrededor de los datos actuales
- Mejorar el almacenamiento de registros y recuperación (cumplimiento con regulaciones)
- Reducir costos alrededor de la plataforma de integración de datos

Tres Escenarios Comunes de Modernización del Almacén de Datos



PureData Systems for Analytics y BigInsights: Fundación de un Almacén de Datos Moderno

InfoSphere BigInsights



PureData System for Analytics



Proporciona...

- Habilidad de fácilmente mover datos entre los dos sistemas
- Habilidad de capturar y procesar nuevos tipos de datos
- Habilidad de consultar todos tipos de datos usando habilidades SQL existentes
- Entorno optimizado de almacén de datos moderno

Service Oriented Finance Quiere Empezar a Crear su Almacén de Datos Moderno

Aquí esta una buen manera de empezar que...

1. Ahorra dinero
2. Libera recursos en su almacén existente
3. Utiliza sus habilidades de SQL existentes



IBM

Paso 1

Mueva datos mas antiguos a BigInsights

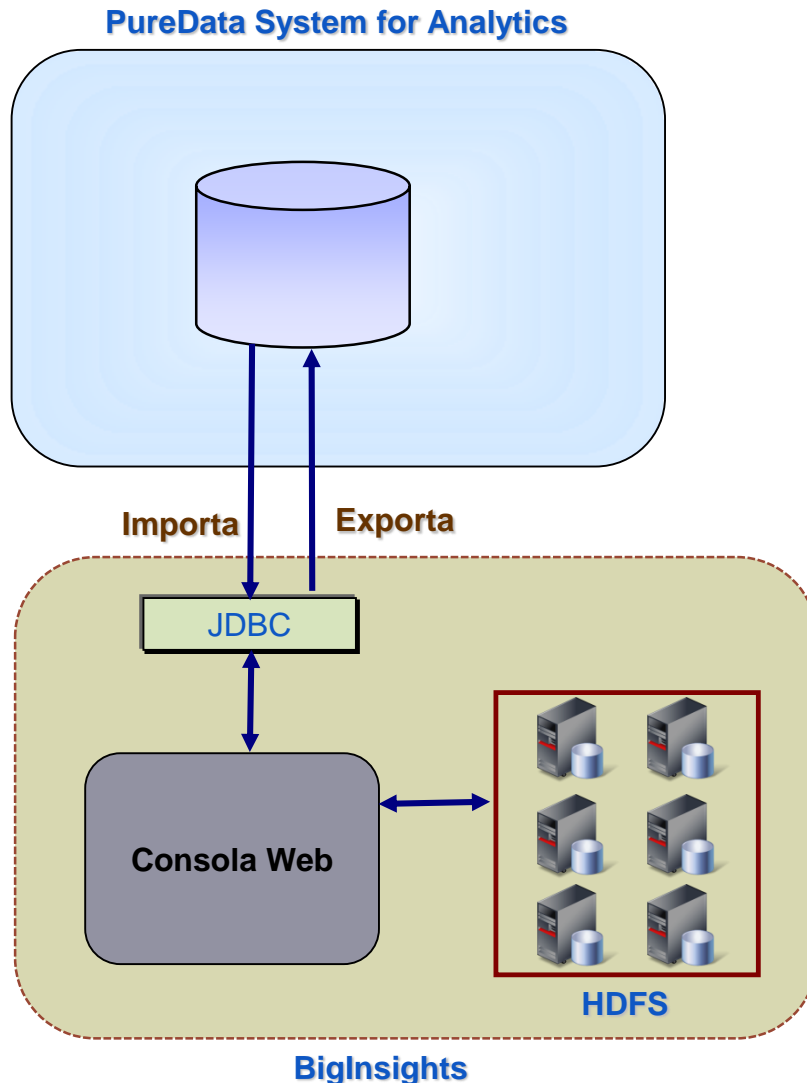
Paso 2

Define el esquema en términos SQL

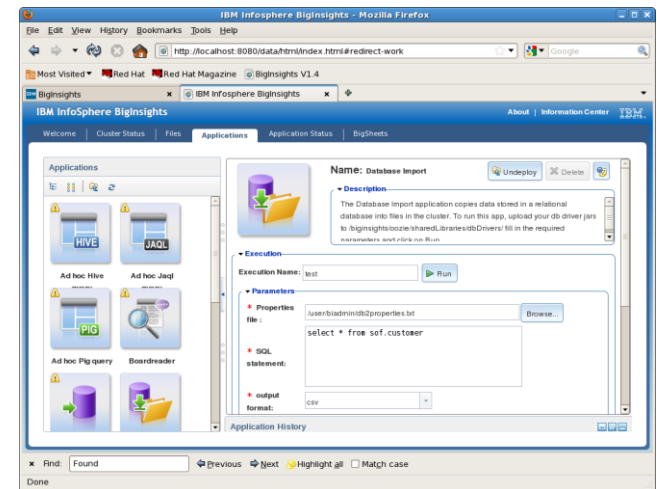
Paso 3

Empiece acezar archivos que se pueden consultar!

Usando BigInsights Console y JDBC – Exporta e Importa

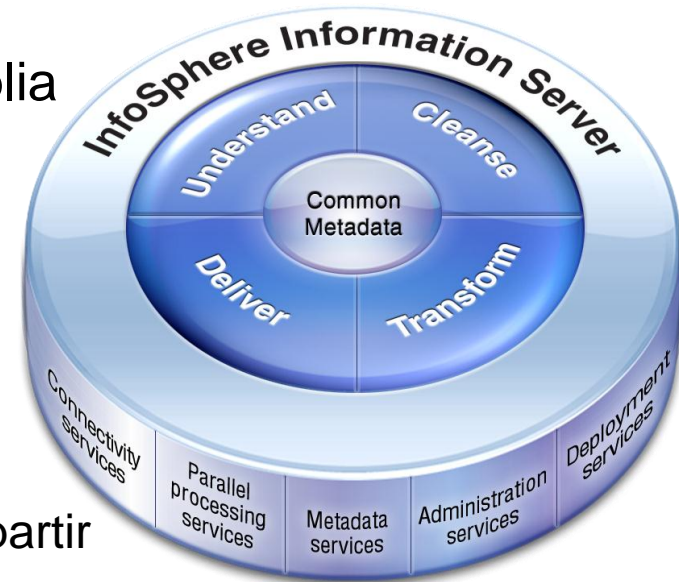


- Parámetros JDBC especificados en la consola de la aplicación
- La importación carga HDFS con datos relacionales – formato es CSV o JSON
- La exportación escribe datos de HDFS en una tabla relacional – Formato HDFS debe ser CSV o JSON



El “Information Server” Proporciona una Solución Integrada Completa de Clase Empresarial

- Se usan las mismas herramientas de diseño graficas independiente de la fuente de los datos
- Herramientas probadas con años de uso y amplia difusión
- Escalabilidad lineal probada
- Acelera la productividad, reduce los costos y riesgos
 - Metadatos compartidos fomentan colaboración
 - Maximiza reutilización – construir una vez y compartir
 - Cienes de componentes pre-construidos
- Se conecta a cualquier fuente de datos
 - Base de datos, aplicaciones, archivos, colas de mensajes, Hadoop, NoSQL

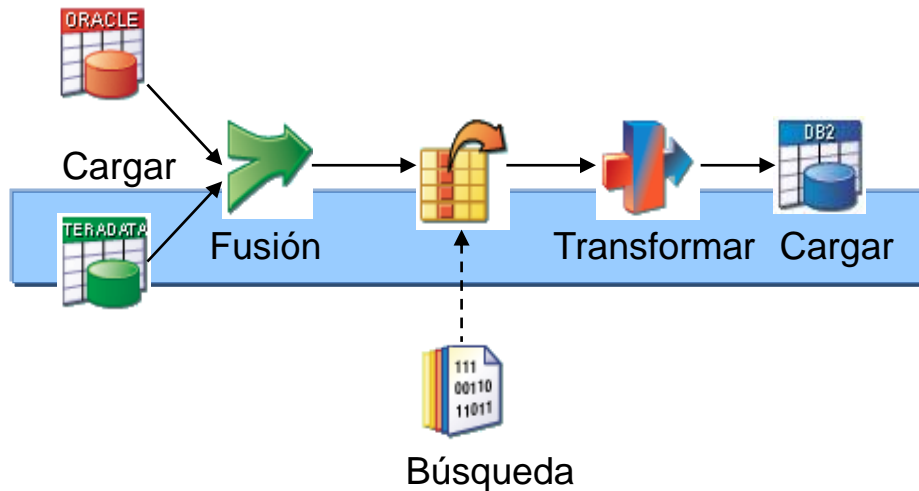


Step 1

Mueva datos mas antiguos a BigInsights

DataStage Hace Fácil La Integración de Todos Las Fuentes de Datos

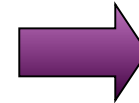
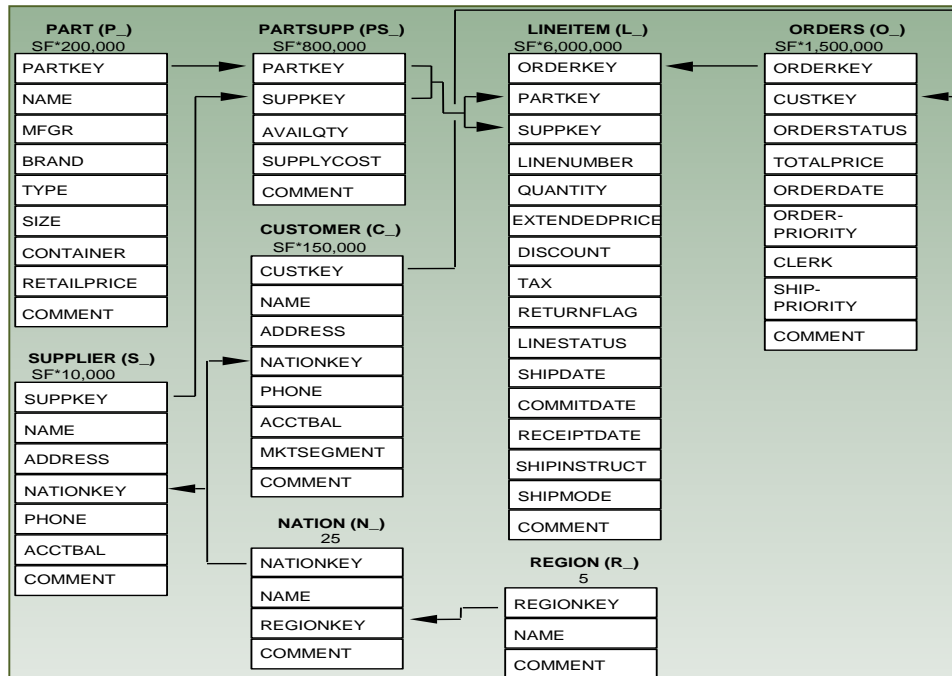
Arrastre, suelte, y configure



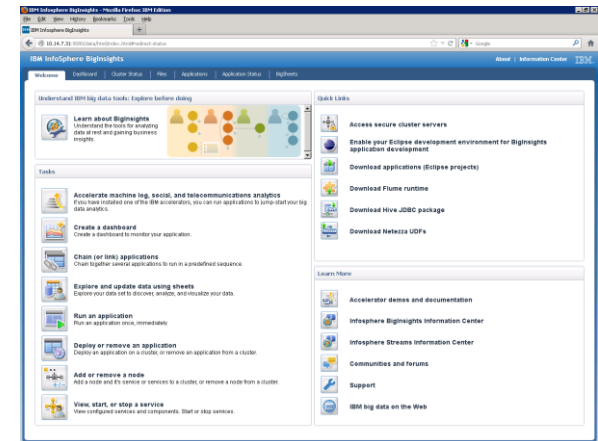
- **Productivo**
 - ▶ Diseño gráfico de los flujos de datos
 - Se centran en el flujo en lugar de la implementación
 - ▶ Numerosos componentes pre-construidos
- **Escalable**
 - ▶ Aprovecha procesamiento paralelo
- **Reduce el riesgo y el costo**
 - ▶ Enfoque modular maximiza la reutilización de los componentes
 - ▶ Metadatos compartidos mejora la colaboración
- **Una herramienta integra todas las fuentes de datos**

Paso 2: Definir Esquema Existente en BigInsights

Esquema Relacional Existente



InfoSphere BigInsights



Paso 2

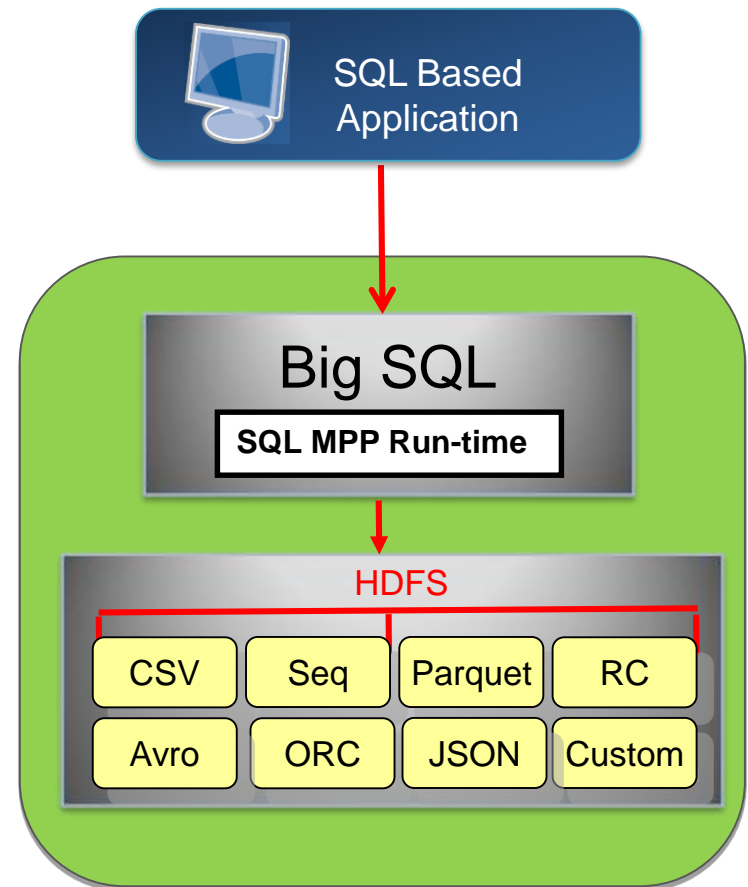
Define el esquema en términos SQL

Big SQL V3.0: Moviendo SQL en Hadoop al Siguiente Nivel

- Masivamente paralelo motor SQL en Hadoop
 - Diseñado desde el principio para baja latencia y alto rendimiento

- Soporte SQL comprensivo
 - El mismo SQL que utiliza en su almacén de datos debe ejecutar con pocas o ninguna modificación
 - Soporte completo para sub-consultas
 - Todas operaciones “JOIN” estándares
 - Procedimientos almacenados / Funciones definidas por usuario

- Soporta todos los formatos de archivos modernos

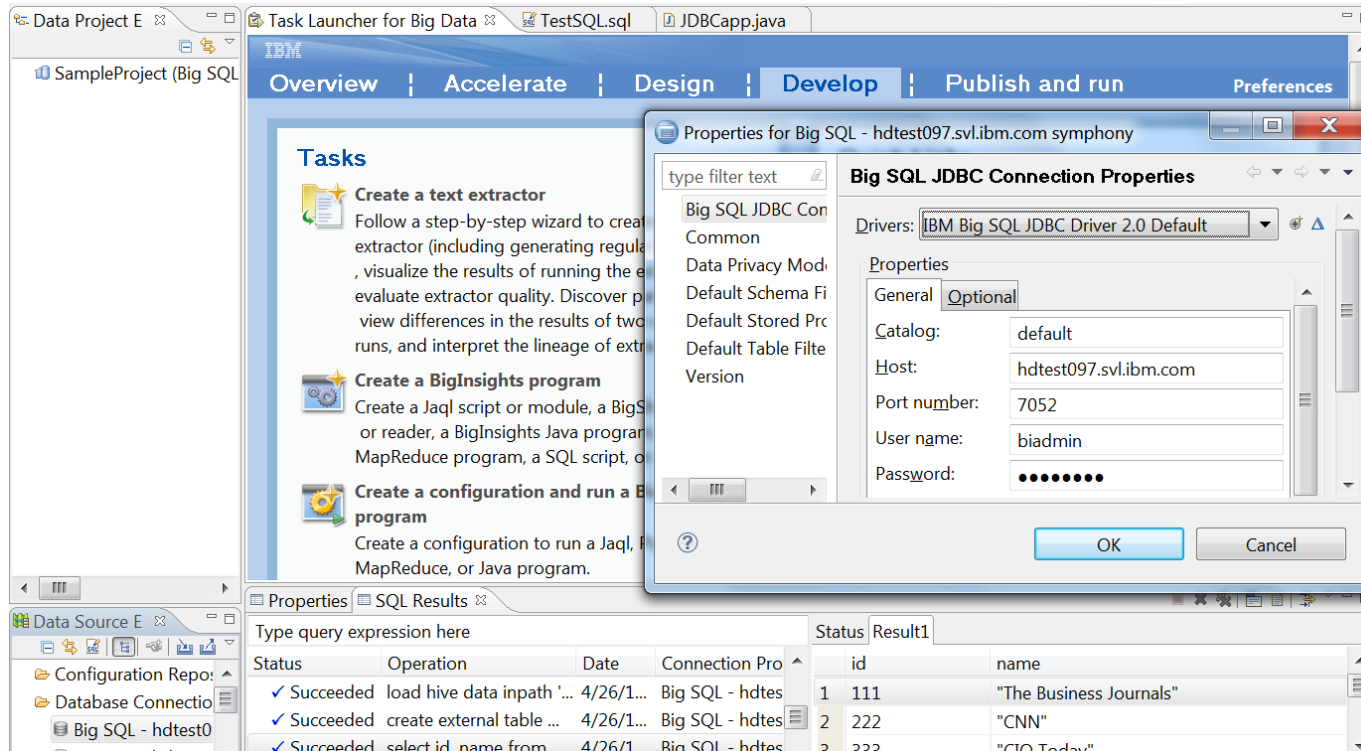


¿Por qué Queremos Utilizar SQL con Datos en Hadoop?

- Programación MapReduce es difícil
 - MapReduce Java API requiere conocimientos de programación
- Hadoop/MapReduce son tecnologías nuevas
 - La cantidad de experiencia es limitada
- Lenguas desconocidas (como Pig) también requieren habilidades especiales
- Compatibilidad con SQL abre los datos a un público mucho más amplio
 - Rampa de abordar fácil para profesionales de Hadoop para SQL
- Soporte para SQL abre los datos a una amplia variedad de herramientas de SQL
 - Cognos, JDBC, ODBC

Opciones de Invocación Proporcionados con BigInsights

- Command-line interface (JSqsh shell)
- Web-based interface (BigInsights web console)
- Eclipse (BigInsights plug-in)



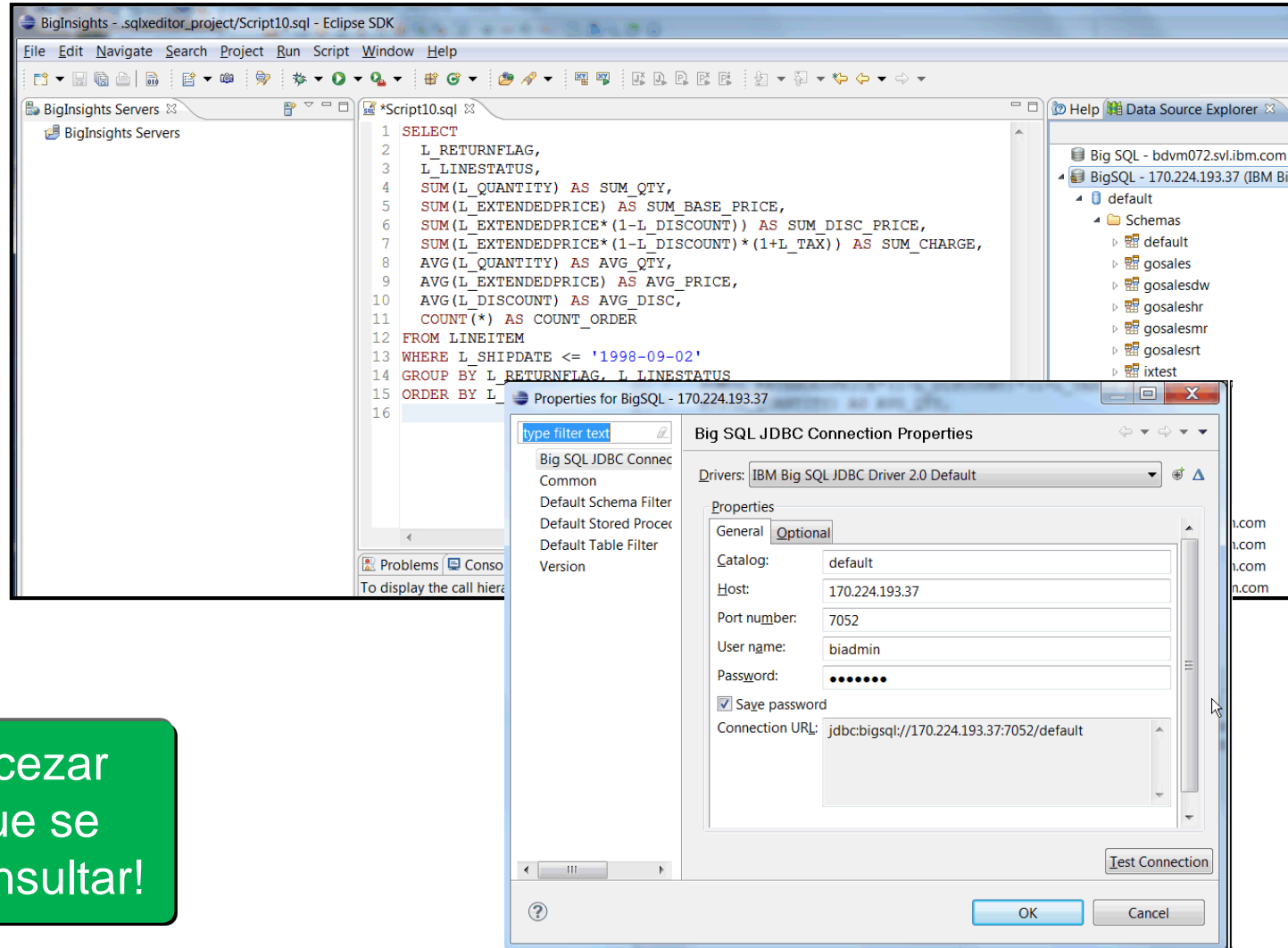
The screenshot displays the Eclipse IDE interface. In the foreground, a dialog box titled "Properties for Big SQL - hdtest097.svl.ibm.com symphony" is open, showing the "Big SQL JDBC Connection Properties" configuration. The "Optional" tab is selected, and the following fields are visible:

- Drivers: IBM Big SQL JDBC Driver 2.0 Default
- General tab selected
- Optional tab selected
- Catalog: default
- Host: hdtest097.svl.ibm.com
- Port number: 7052
- User name: biadmin
- Password: [masked]

In the background, the Eclipse IDE shows a "Tasks" panel with three tasks: "Create a text extractor", "Create a BigInsights program", and "Create a configuration and run a BigInsights program". Below the dialog, the "SQL Results" table is visible, showing the following data:

Status	Operation	Date	Connection Pro	id	name
✓ Succeeded	load hive data inpath '...	4/26/1...	Big SQL - hdt...	1 111	"The Business Journals"
✓ Succeeded	create external table ...	4/26/1...	Big SQL - hdt...	2 222	"CNN"
✓ Succeeded	select id, name from ...	4/26/1...	Big SQL - hdt...	3 333	"CIO Today"

DEMO: Usando Herramientas Eclipse Para el Acceso de Archivos que se Pueden Consultar



The screenshot shows the Eclipse IDE interface. The main editor displays a SQL script with the following content:

```

1 SELECT
2   L_RETURNFLAG,
3   L_LINESTATUS,
4   SUM(L_QUANTITY) AS SUM_QTY,
5   SUM(L_EXTENDEDPRI) AS SUM_BASE_PRICE,
6   SUM(L_EXTENDEDPRI*(1-L_DISCOUNT)) AS SUM_DISC_PRICE,
7   SUM(L_EXTENDEDPRI*(1-L_DISCOUNT)*(1+L_TAX)) AS SUM_CHARGE,
8   AVG(L_QUANTITY) AS AVG_QTY,
9   AVG(L_EXTENDEDPRI) AS AVG_PRICE,
10  AVG(L_DISCOUNT) AS AVG_DISC,
11  COUNT(*) AS COUNT_ORDER
12 FROM LINEITEM
13 WHERE L_SHIPDATE <= '1998-09-02'
14 GROUP BY L_RETURNFLAG, L_LINESTATUS
15 ORDER BY L_
16

```

Overlaid on the IDE is the "Big SQL JDBC Connection Properties" dialog box. The "Drivers" dropdown is set to "IBM Big SQL JDBC Driver 2.0 Default". The "Optional" tab is selected, showing the following configuration:

- Drivers: IBM Big SQL JDBC Driver 2.0 Default
- Properties:
 - General
 - Optional
- Catalog: default
- Host: 170.224.193.37
- Port number: 7052
- User name: biadmin
- Password: ••••••
- Save password
- Connection URL: jdbc:bigsql://170.224.193.37:7052/default

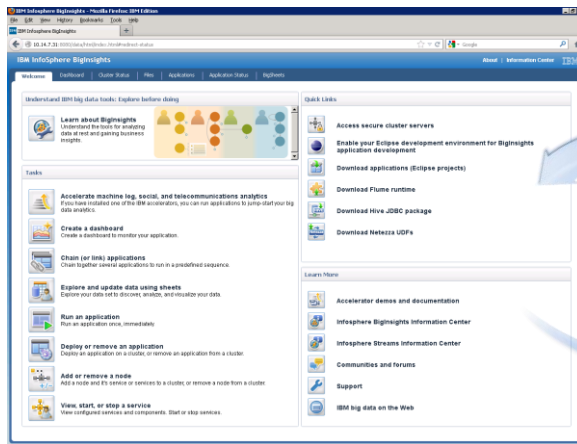
Buttons for "Test Connection", "OK", and "Cancel" are visible at the bottom of the dialog.

Step 3

Empiece a cezar archivos que se pueden consultar!

Service Oriented Finance ha Creado un Almacén de Datos Moderno

InfoSphere BigInsights



PureData System for Analytics



Paso 1

Mueva datos mas antiguos a BigInsights

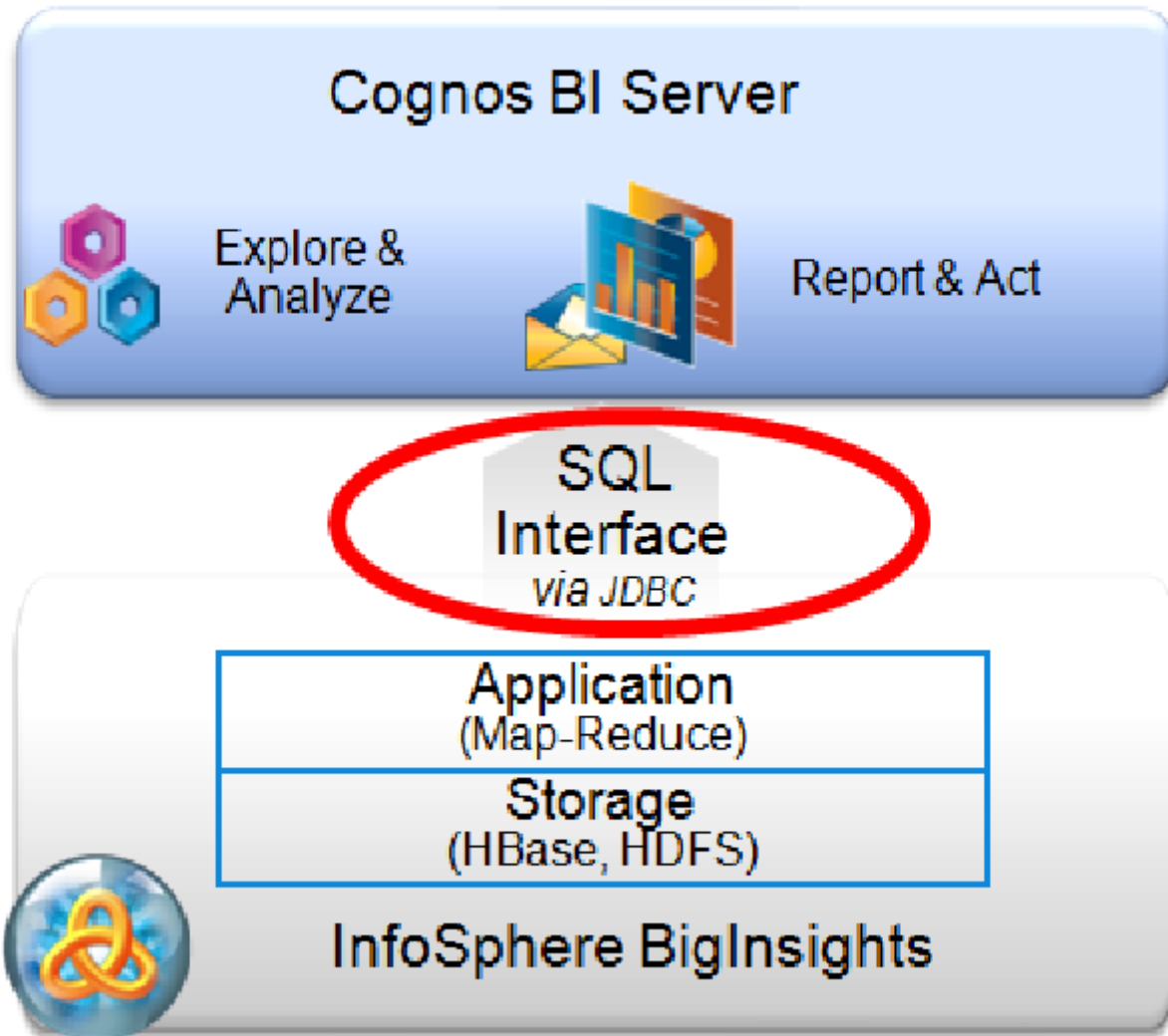
Paso 2

Define el esquema en términos SQL

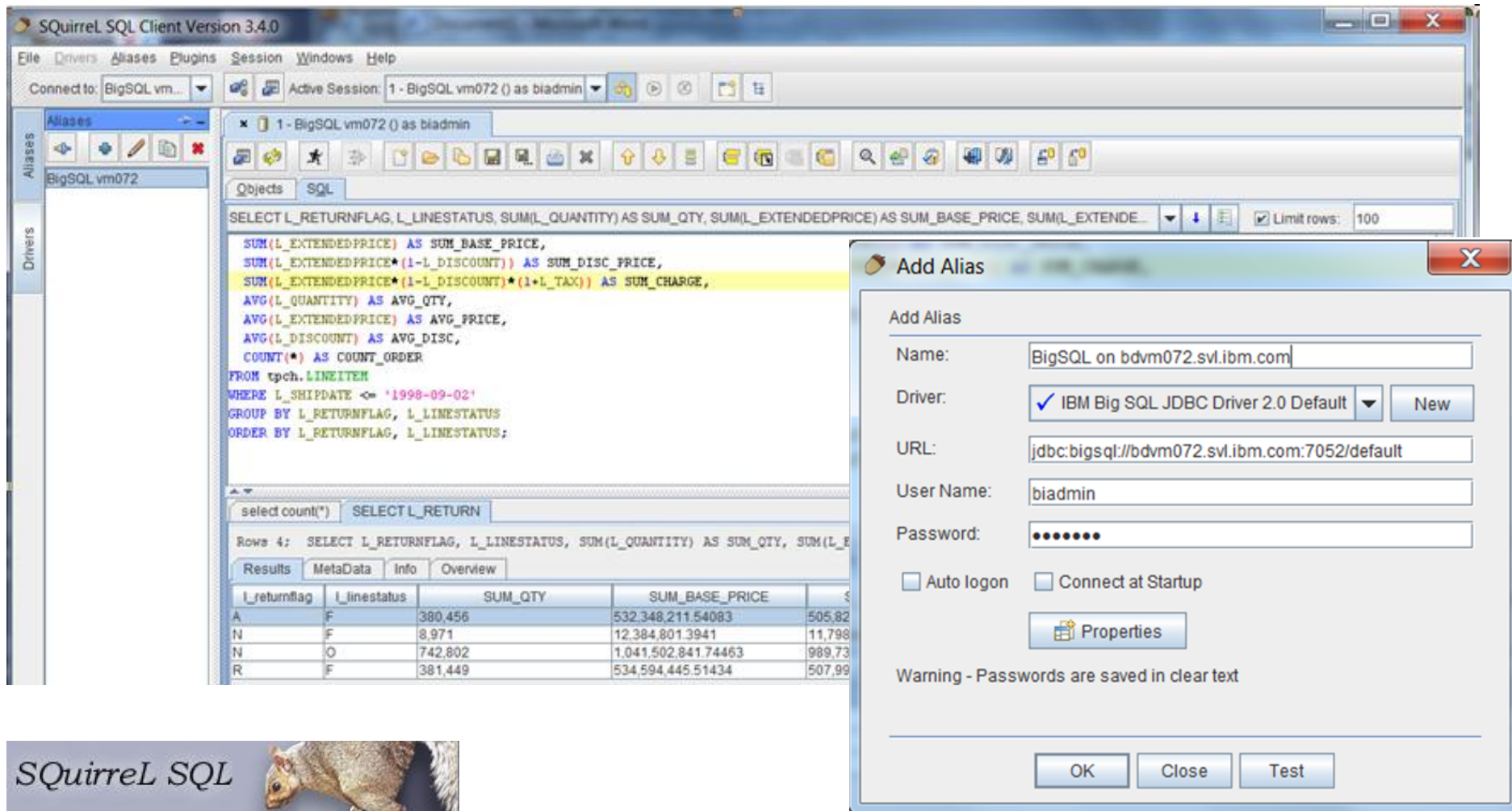
Paso 3

Empiece a cezar archivos que se pueden consultar!

Inteligencia Empresarial Cognos



Usando Existente Herramientas Estándares Como Squirrel SQL



The screenshot shows the Squirrel SQL Client interface. The main window displays a SQL query and its results. An 'Add Alias' dialog box is open in the foreground, showing the configuration for a new alias.

SQL Query:

```
SELECT L_RETURNFLAG, L_LINESTATUS, SUM(L_QUANTITY) AS SUM_QTY, SUM(L_EXTENDEDPRICE) AS SUM_BASE_PRICE, SUM(L_EXTENDEDPRICE*(1-L_DISCOUNT)) AS SUM_DISC_PRICE, SUM(L_EXTENDEDPRICE*(1-L_DISCOUNT)*(1+L_TAX)) AS SUM_CHARGE, AVG(L_QUANTITY) AS AVG_QTY, AVG(L_EXTENDEDPRICE) AS AVG_PRICE, AVG(L_DISCOUNT) AS AVG_DISC, COUNT(*) AS COUNT_ORDER
FROM tpch.LINEITEM
WHERE L_SHIPDATE <= '1998-09-02'
GROUP BY L_RETURNFLAG, L_LINESTATUS
ORDER BY L_RETURNFLAG, L_LINESTATUS;
```

Add Alias Dialog Box:

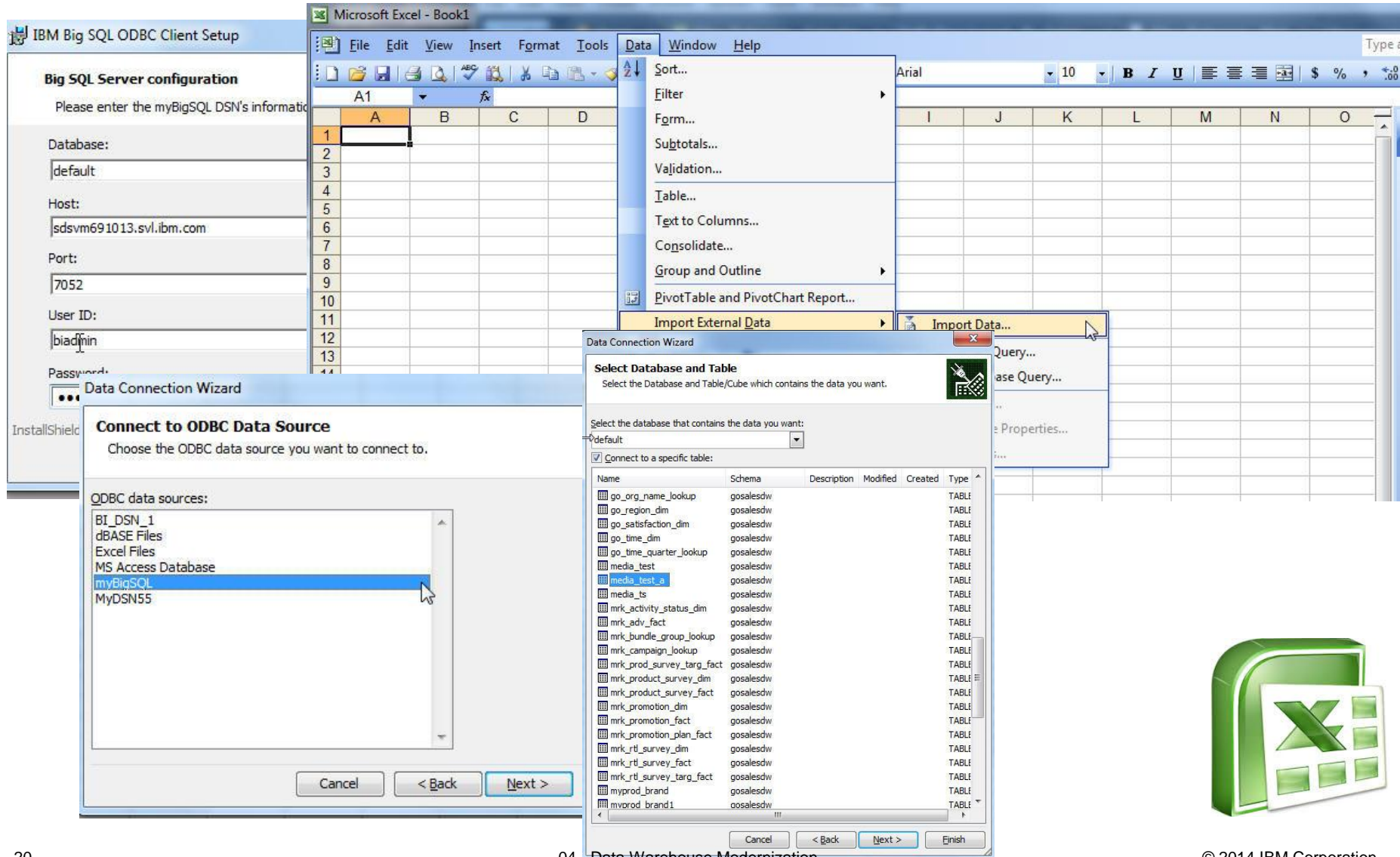
- Name: BigSQL on bdvm072.svl.ibm.com
- Driver: IBM Big SQL JDBC Driver 2.0 Default New
- URL: jdbc:bigsql://bdvm072.svl.ibm.com:7052/default
- User Name: biadmin
- Password: [REDACTED]
- Auto logon Connect at Startup
- Properties

Results Table:

L_returnflag	L_linestatus	SUM_QTY	SUM_BASE_PRICE	SUM_DISC_PRICE	SUM_CHARGE
A	F	380,456	532,348,211.54083	505,82	
N	F	8,971	12,384,801.3941	11,798	
N	O	742,802	1,041,502,841.74463	989,73	
R	F	381,449	534,594,445.51434	507,99	



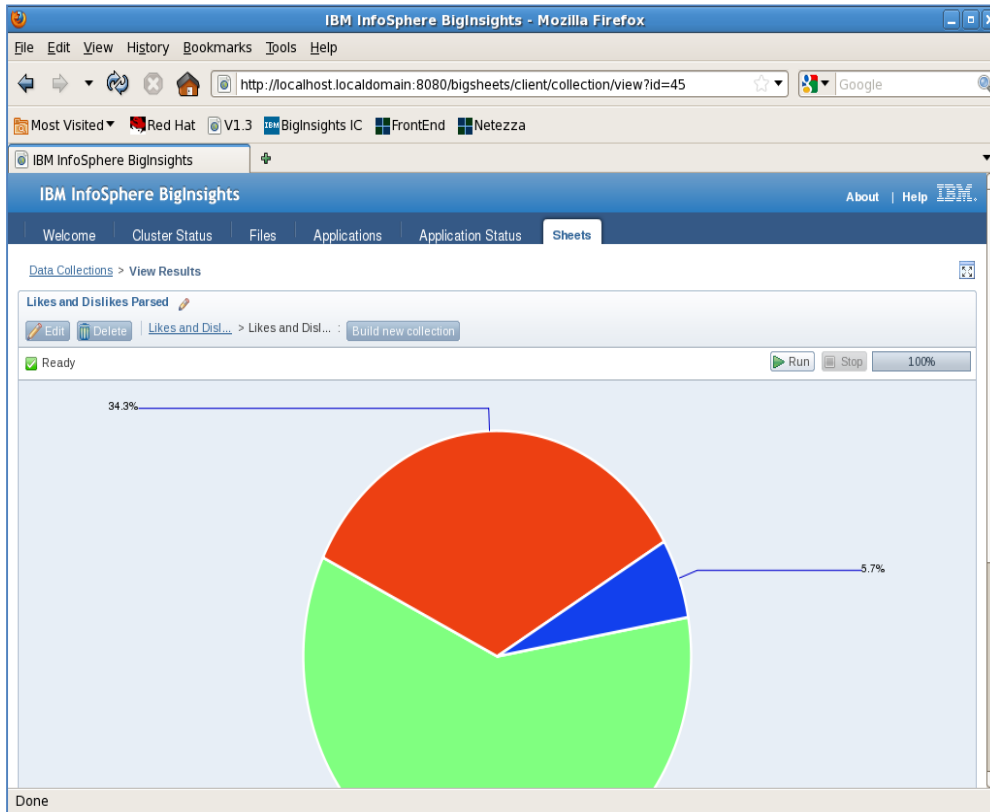
MS Excel: Integración Big SQL a Través de ODBC



The screenshot illustrates the integration of Big SQL into Microsoft Excel via ODBC. The main window shows Microsoft Excel with the 'Data' menu open, highlighting 'Import External Data'. A 'Data Connection Wizard' dialog is active, showing the 'Select Database and Table' step. The 'Connect to ODBC Data Source' dialog is also open, showing a list of ODBC data sources with 'myBigSQL' selected. The 'Data Connection Wizard' shows a list of tables from the 'default' database, including 'go_org_name_lookup', 'go_region_dim', 'go_satisfaction_dim', 'go_time_dim', 'go_time_quarter_lookup', 'media_test', 'media_test_a', 'media_ts', 'mrk_activity_status_dim', 'mrk_adv_fact', 'mrk_bundle_group_lookup', 'mrk_campaign_lookup', 'mrk_prod_survey_targ_fact', 'mrk_product_survey_dim', 'mrk_product_survey_fact', 'mrk_promotion_dim', 'mrk_promotion_fact', 'mrk_promotion_plan_fact', 'mrk_rtl_survey_dim', 'mrk_rtl_survey_fact', 'mrk_rtl_survey_targ_fact', 'myprod_brand', and 'myprod_brand1'.



Datos Disponibles para Todas las Herramientas de BigInsights



- BigSheets spreadsheet y visualización
- Acelerador de datos de maquina
- Acelerador de datos de medios sociales
- Analítica de texto avanzado
- Lenguaje de consulta JAQL
- Consola de Herramientas
- Big R –Integración con lenguaje R

BigInsights con Big SQL en POWER Proporciona Consultas Big Data Mas Rápido Que Hive en x86

Consulta de Carga de Trabajo de Inteligencia Empresarial Moderno

BigInsights v3.0 en 8 S822L nodos

S822L, 3.3 GHz
24 núcleos
256 GB Memoria
RHEL 6.5



Tiempo para completar las consultas en 7 corrientes concurrentes

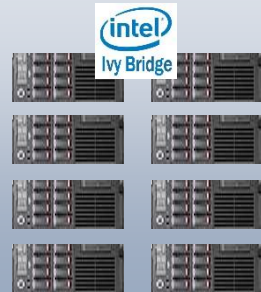
**8 horas
40 min**

**11.2x
Mas rápido**

**4.4x
Menor precio de
rendimiento**

Hive v0.12 en 8 Ivy Bridge EP nodos

Intel Ivy Bridge EP
2.7GHz
24 núcleos
256 GB Memoria
RHEL 6.4



Tiempo para completar las consultas en 7 corrientes concurrentes

**97 horas
27 min**

This is an IBM internal study of a benchmark-inspired workload in a controlled laboratory environment. The IBM system under test consists of eight Power S822L data nodes (24 cores each @ 3.3 GHz) with 256 GB RAM running RHEL 6.5. The Intel-based system under test consists of eight Ivy Bridge EP data nodes (24 cores each @ 2.7 GHz) with 256 GB RAM running RHEL 6.4. Both systems each used a 10TB scale factor for the TCP-DS inspired test. Both systems were driven by a single management system (Power S822L with 24 cores, 256 GB RAM, RHEL 6.5 and Ivy Bridge EP - 24 cores, 256 GB RAM, RHEL 6.4) that were immaterial to the performance of the systems under test. Tests measured total elapsed time to complete an identical subset of 42 of the 99 queries in the TPC-DS-inspired workload. Test executed 7 simultaneous query streams. Measured time was the maximum of the 7 streams across 2 back-to-back runs. Customer applications, differences in the systems deployed and other system variations or testing conditions may produce different results. Cost analysis based on 3 year total cost of acquisition of hardware, software and support services over a 3 year period. Prices for both IBM and competitor systems based on US list prices valid as of July 2014.

Como Empezar con la Modernización del Almacén de Datos

Consiga educación:

- Descargue BigInsights Quick Start ibm.com/infosphere/quickstart
- Descargue el white paper: [Data Warehouse Augmentation: the Queryable Data Store](#)
- Visite el sitio: [Data Warehouse Augmentation website](#)
- Visite el sitio: IBMBigDataHub.com para acaezar podcasts y video de este caso
- Descargue la investigación: [Capitalizing on Big Data](#)



The collage features several key educational resources:

- The Big Data Hub**: A website with a navigation menu (Home, Blogs, Videos & Podcasts, Resources, Events) and a featured article titled "Analytics: The real-world use of big data".
- White Paper**: "Analytics: The real-world use of big data" by the IBM Institute for Business Value, published in collaboration with the Said Business School at the University of Oxford.
- Big Data University**: A learning platform with sections for Home, Learn, Download, Resources, and Jobs. It promotes "Easy and Affordable Learning Hadoop and other Big Data technologies" and offers a "FREE! Hadoop Fundamentals" course.
- Books**: Several covers are shown, including "Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data" and "Harness the Power of Big Data: The IBM Big Data Platform".
- Testimonials**: A section titled "Student Testimonials" from Balázs (USA) praising the training material for being short, easy to digest, and supported by transcripts and exercises.

Agenda de Hoy

Time	Topic
09:00 – 09:15 AM	Introducción: Lo Que la Analítica Big Data Puede Hacer Para Su Negocio
09:15 – 10:00 AM	Domine los Fundamentos: Analizando datos estructurados con sistemas PureData System for Analytics
10:00 – 10:15 AM	Break
10:15 – 11:00 AM	La Analítica de Datos no Estructurados: Análisis Big Data con Hadoop
11:00 – 11:45 AM	Amplía tu Estrategia de Análisis: Modernización de Almacén de Datos
11:45 – 12:00 PM	Resumen y Acción