

**SAS FORUM ARGENTINA 2016**

**SAS EN HADOOP**  
**CREANDO UN LABORATORIO**  
**BIG DATA / HADOOP / SAS**



**Sergio Uassouf**  
Líder de Práctica de  
Gestión de Información e Infraestructura

## BIG DATA = HADOOP + SAS

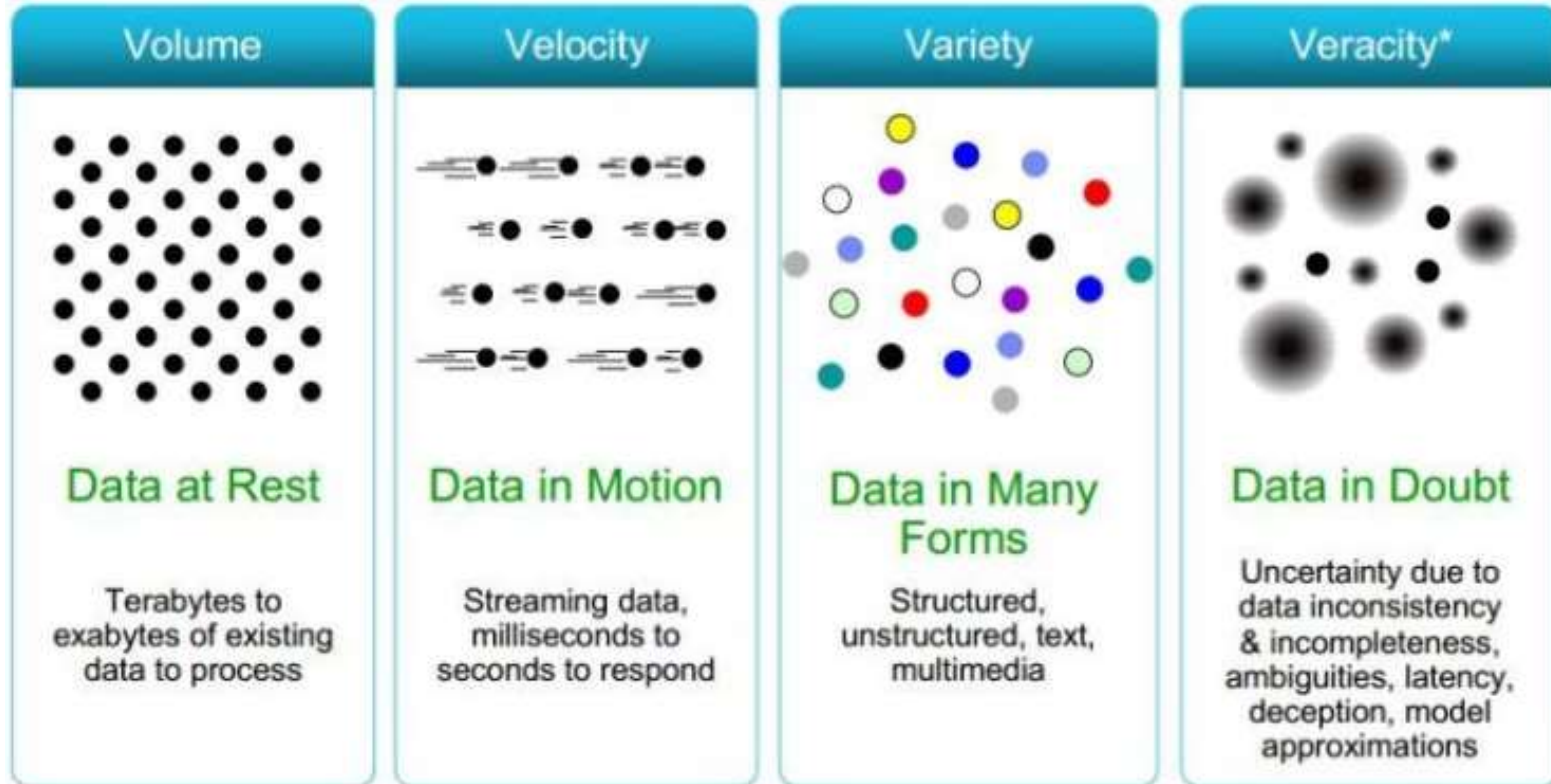
## ¿QUE GRANDE TIENE QUE SER LO GRANDE PARA SER GRANDE?

- ❑ 1024 bytes = 1 Kilobyte
  - ❑ 1024 KB = 1 Megabyte
  - ❑ 1024 MB = 1 Gigabyte
  - ❑ 1024 GB = 1 Terabyte
  - ❑ 1024 TB = 1 Petabyte
  - ❑ 1024 PB = 1 Exabyte
  - ❑ 1024 EB = 1 Zettabyte
  - ❑ 1024 ZB = 1 Yottabyte
  - ❑ 1 YB =  $10^{24}$  bytes
- ❑ Una transacción de cajero automático o home banking = 200 bytes.
  - ❑ 1 TB = 5.000.000.000 de transacciones
  - ❑ Red Link + Banelco ejecutan aprox. 20.000.000 de transacciones por día
  - ❑ 1 TB almacena 250 días de transacciones de Link y Banelco

# Entonces... ¿Qué es Big Data?

# BIG DATA = HADOOP + SAS

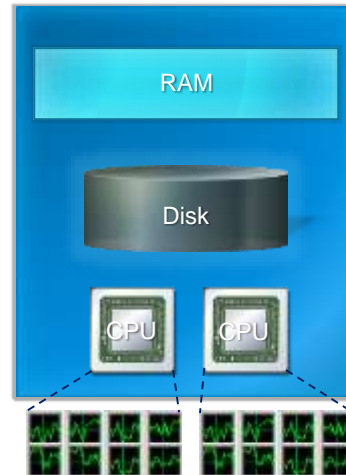
## DEFINICION DE MARKETING



# BIG DATA = HADOOP + SAS

!!! ESTO ES UN PARADIGMA !!!

- Desde los inicios de la informática un computador, ya sea personal o empresarial está compuesto de 3 componentes principales.



MEMORIA

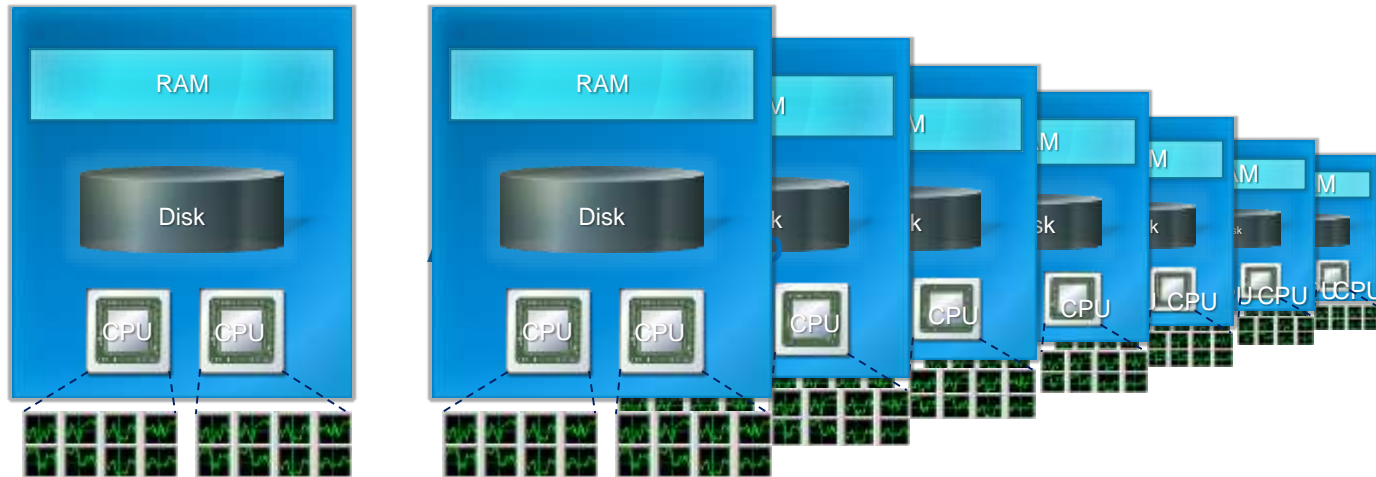
UNIDADES DE  
ALMACENAMIENTO

UNIDADES DE  
PROCESAMIENTO

# BIG DATA = HADOOP + SAS

**ESTO ES BIG DATA**

- Desde los inicios de la informática un computador, ya sea personal o empresarial está compuesto de 3 componentes principales.



**BIG DATA = PROCESAMIENTO MASIVAMENTE PARALELO**

- Si puede almacenar mucha más información a un costo mucho menor...
- Y puede procesarla en un tiempo mucho menor.
- Entonces no necesita armar modelos tomando sólo un subconjunto de los datos...
- Y puede hacer todas las iteraciones que necesite.
- **Entonces puede almacenar y procesar la información que antes no podía**

**BIG DATA =  
HADOOP + SAS**

**ALMACENAR Y ANALIZAR  
GRANDES VOLUMENES DE INFORMACIÓN  
A BAJO COSTO**

**TODOS LOS  
CALL DETAIL  
RECORDS**

**TODAS LAS  
TRANSACCIONES**

**TODAS LAS  
SECUENCIAS DE  
SITIOS WEB**

**TODAS LAS  
CONVERSACIONES  
DE LOS CALL  
CENTERS**

**Y ANALIZARLOS  
EN SU TOTALIDAD...**

**EJECUTANDO  
TODAS LAS  
ITERACIONES QUE  
NECESITE...**

**A MUY BAJO  
COSTO RELATIVO**



# BIG DATA = HADOOP + SAS

INTERFAZ DE  
USUARIO

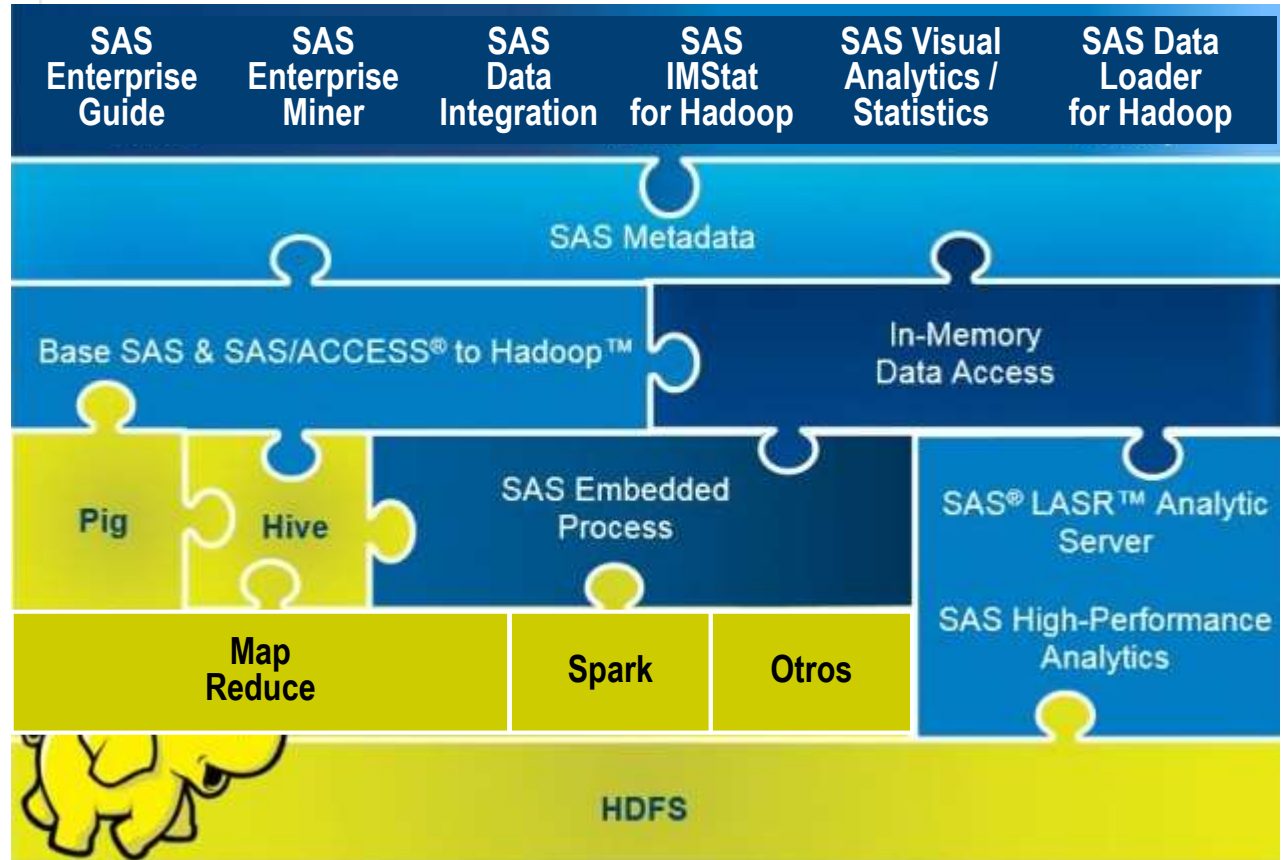
METADATOS

ACCESO A  
DATOS

PROCESAMIENTO  
DE DATOS

FILE SYSTEM

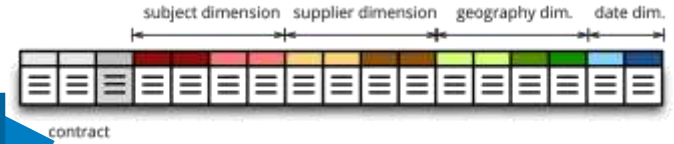
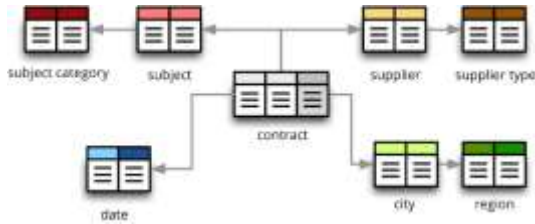
## COMPONENTES DEL "ECOSISTEMA"





# BIG DATA = HADOOP + SAS

## PREPARACIÓN Y ANÁLISIS DE DATOS

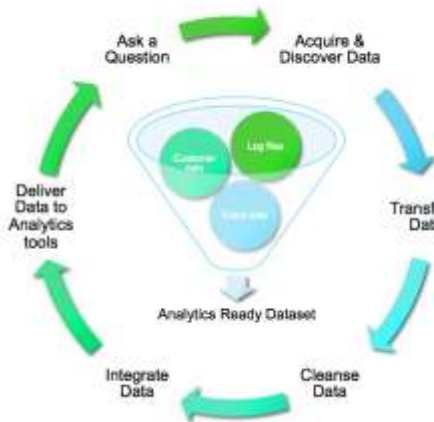


# BIG DATA = HADOOP + SAS

# PREPARACIÓN DE DATOS "IN-DATABASE" ESTADÍSTICA Y MINERÍA DE DATOS "IN-MEMORY"

Data Preparation Process

Analytics & Visualization



## SAS Data Loader for Hadoop

SAS Data Loader

What directive do you want to perform?

- Search Procedures**: Open up existing SAS procedures in your environment.
- Run Scripts**: Show the status of scripts and process details procedures.
- Change on Job Data in Hadoop**: Change a data set in Hadoop.
- Run a SAS Program**: Run a SAS program in your environment.
- Transform Data in Hadoop**: Transform data in Hadoop.
- Transfer Data in Hadoop**: Transfer data in Hadoop.
- Copy Data to Hadoop**: Copy data to Hadoop.
- Export Data to SAS**: Export data to SAS.
- Create Jobs in Hadoop**: Create jobs in Hadoop.

## Visual Analytics / Statistics



## Applications Run Natively **IN** Hadoop

**DATA LOADER FOR HADOOP**  
(EMBEDDED PROCESS - LENGUAJE DS2)

**SAS VISUAL ANALYTICS / STATISTICS**  
(LASR SERVER - IMSTAT)

**YARN** (Cluster Resource Management)

**HDFS2** (Redundant, Reliable Storage)



# HADOOP PROCESAMIENTO "IN-DATABASE"

- ❑ Embedded Process: Traduce código de alto nivel a Map-Reduce.

```
proc ds2 ;  
/* thread ~ equiv to a mapper */  
  thread map_program;  
  method run(); set dbmslib.intab;  
  /* program steps */  
end; endthread;  
/* program steps */  
data hdf.data;  
dcl thread map_program map_pgm; method  
run();  
set from map_pgm threads=N;  
/* reduce steps */ end; enddata;  
run; quit;
```



1. Integración
2. Preparación
3. Calidad
4. Scoring



# DATA LOADER FOR HADOOP

## IDEAL PARA PREPARACIÓN DE TABLAS ANALITICAS

### CODIGO ALTO NIVEL EN HADOOP

```
proc ds2 ;  
/* thread ~ equiv to a mapper */  
  thread map_program;  
  method run(); set dbmslib.intab;  
  /* program statements */  
end; endthread; run;  
/* program wrapper */  
data hdf.data_reduced;  
dcl thread map_program map_pgm; method  
run();  
set from map_pgm threads=N;  
/* reduce steps */ end; enddata;  
run; quit;
```

### DATA LOADER FOR HADOOP



**Saved Directives**  
Open a previously created directive to run, view or edit.



**Run Status**  
Show the status of current and previous directive executions



**Query or Join Data in Hadoop**  
Query a table, or join data from multiple tables



**Sort and De-Duplicate Data in Hadoop**  
Query, sort, or de-duplicate the data in an existing Hadoop table



**Run a SAS Program**  
Run in-database data quality SAS programs



**Transform Data in Hadoop**  
Transform data from a Hadoop table



**Transpose Data in Hadoop**  
Transpose data from a Hadoop table



**Copy Data from Hadoop**  
Copy Data from Hadoop into a database



**Copy Data to Hadoop**  
Copy data from a database into Hadoop



**Load Data to LASR**  
Copy data from a source and load it into LASR. Existing data in the target table will be replaced.



**Cleanse Data in Hadoop**  
Cleanse data in Hadoop by performing data quality transforms



**Profile Data**  
Generate a profile report of the data in a table



**Saved Profile Reports**  
Explore previously generated profile reports

# DATA LOADER FOR HADOOP

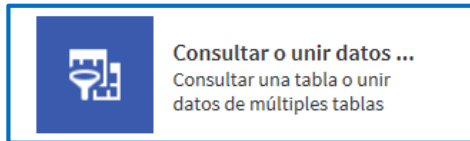
## MENÚ PRINCIPAL DE DIRECTIVAS



**Directivas guardadas**  
Abrir una directiva ya creada para ejecutar, ver o editar



**Estado de ejecución**  
Mostrar el estado de las ejecuciones de directivas anteriores y actuales



**Consultar o unir datos ...**  
Consultar una tabla o unir datos de múltiples tablas



**Ordenar y deduplicar d...**  
Consultar, ordenar o deduplicar datos de una tabla Hadoop existente



**Ejecutar programas SAS**  
Ejecutar programas SAS de calidad de datos en la base de datos



**Transformar datos en H...**  
Transformar datos de una tabla Hadoop



**Transponer datos en H...**  
Transponer datos de una tabla Hadoop



**Copiar datos desde Ha...**  
Copiar datos desde Hadoop en una base de datos



**Copiar datos en Hadoop**  
Copiar datos de una base de datos en Hadoop



**Cargar datos en LASR**  
Copiar datos de una fuente de datos y cargarlos en LASR. Se reemplazarán los datos exis...



**Limpiar datos en Hadoop**  
Limpiar datos en Hadoop realizando transformaciones de calidad de datos



**Datos del perfil**  
Generar un informe de perfil de los datos en una tabla



**Informes de perfil guar...**  
Explorar informes de perfil previamente generados

# DATA LOADER FOR HADOOP

## QUERIES Y JOINS



Query or Join Data in Hadoop  
Query a table, or join data from multiple tables

**JOIN** *Inner Join: cars.make = sample\_07.code*

*Choose a table to query, or multiple tables to join and the columns to join on*

Base table:  ...

Join:  Inner Join  Left Join  Right Join  Full Join

Join on:  Inner Join  Left Join  Right Join  Full Join

...  =  ...

Next



# DATA LOADER FOR HADOOP

## MENÚ PRINCIPAL DE DIRECTIVAS



**Directivas guardadas**  
Abrir una directiva ya creada para ejecutar, ver o editar



**Estado de ejecución**  
Mostrar el estado de las ejecuciones de directivas anteriores y actuales



**Consultar o unir datos ...**  
Consultar una tabla o unir datos de múltiples tablas



**Ordenar y deduplicar d...**  
Consultar, ordenar o deduplicar datos de una tabla Hadoop existente



**Ejecutar programas SAS**  
Ejecutar programas SAS de calidad de datos en la base de datos



**Transformar datos en H...**  
Transformar datos de una tabla Hadoop



**Transponer datos en H...**  
Transponer datos de una tabla Hadoop



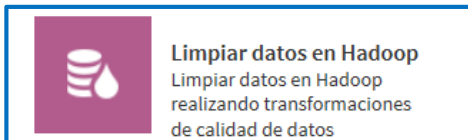
**Copiar datos desde Ha...**  
Copiar datos desde Hadoop en una base de datos



**Copiar datos en Hadoop**  
Copiar datos de una base de datos en Hadoop



**Cargar datos en LASR**  
Copiar datos de una fuente de datos y cargarlos en LASR. Se reemplazarán los datos exis...



**Limpiar datos en Hadoop**  
Limpiar datos en Hadoop realizando transformaciones de calidad de datos



**Datos del perfil**  
Generar un informe de perfil de los datos en una tabla



**Informes de perfil guar...**  
Explorar informes de perfil previamente generados

# DATA LOADER FOR HADOOP

## MENÚ DATA CLEANSING



### Limpiar datos en Hadoop

Limpiar datos en Hadoop realizando transformaciones de calidad de datos



### Filter Data

Select the rows of data to include



### Generate Match Codes

Create match codes for selected values in the table



### Identification Analysis

Identify the semantic data type of text in selected columns



### Manage Columns

Select the columns to include



### Parse Data

Select the column, Definition, and Token you want to apply, and enter a name for the new column



### Standardize Data

Apply data standards to selected columns



### Summarize Rows

Create a new row with data summarized in selected columns

## Applications Run Natively **IN** Hadoop

SAS DATA LOADER FOR HADOOP  
(EMBEDDED PROCESS - SAS LENGUAJE DS2)

SAS VISUAL ANALYTICS / STATISTICS  
(LASR SERVER - IMSTAT)

**YARN** (Cluster Resource Management)

**HDFS2** (Redundant, Reliable Storage)



# SAS LASR ANALYTICS SERVER

## Data Manipulation

- SAS Data Step
- BALANCE
- COLUMNINFO
- COMPUTE
- DELETEROWS
- DISTINCT
- DROPTABLE
- FETCH
- GROUPBY
- PARTITION
- PROMOTE
- PURGETEMPFILES
- SET
- TABLE
- UPDATE

## Data Exploration/ Visualization

- BOXPLOT
- CORR
- CROSSTAB
- CONTOURPLOT
- DISTRIBUTIONINFO
- FREQUENCY
- HISTOGRAM
- KDE
- REPLAY
- SUMMARY

## Predictive Modeling

- DECISIONTREE
- FORECAST
- GENMODEL
- GLM
- RANDOMWOODS
- ASSESSMENT

## Descriptive Modeling

- CLUSTER
- CLUSTER TF-IDF
- ASSOCIATIONS
- SVD

## Recommender

- CLUSTER
- KNN
- ASSOCIATIONS
- SVD

## Text Analytics

- PARSING
- SVD

## Miscellaneous

- EXTERNAL (C API)
- FREE
- SAVE
- STORE

## Deployment

- SCORE

**Data  
Manipulation**

**Exploration/  
Visualization**

**Modeling**

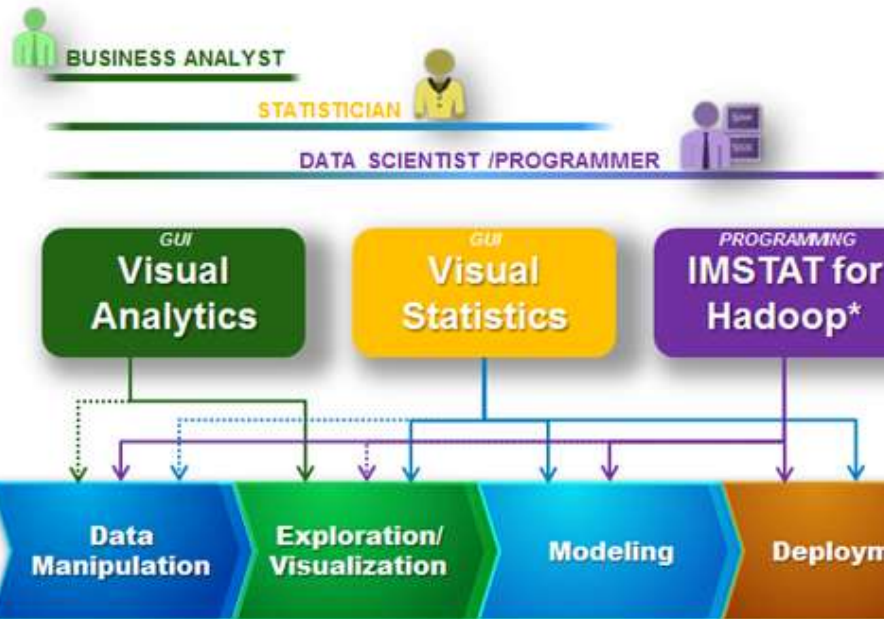
**Deployment**

## IN-MEMORY STATISTICS

```
data example.iris;
  set sashelp.iris;
run;
proc imstat data=example.iris;
  corr;
quit;
```

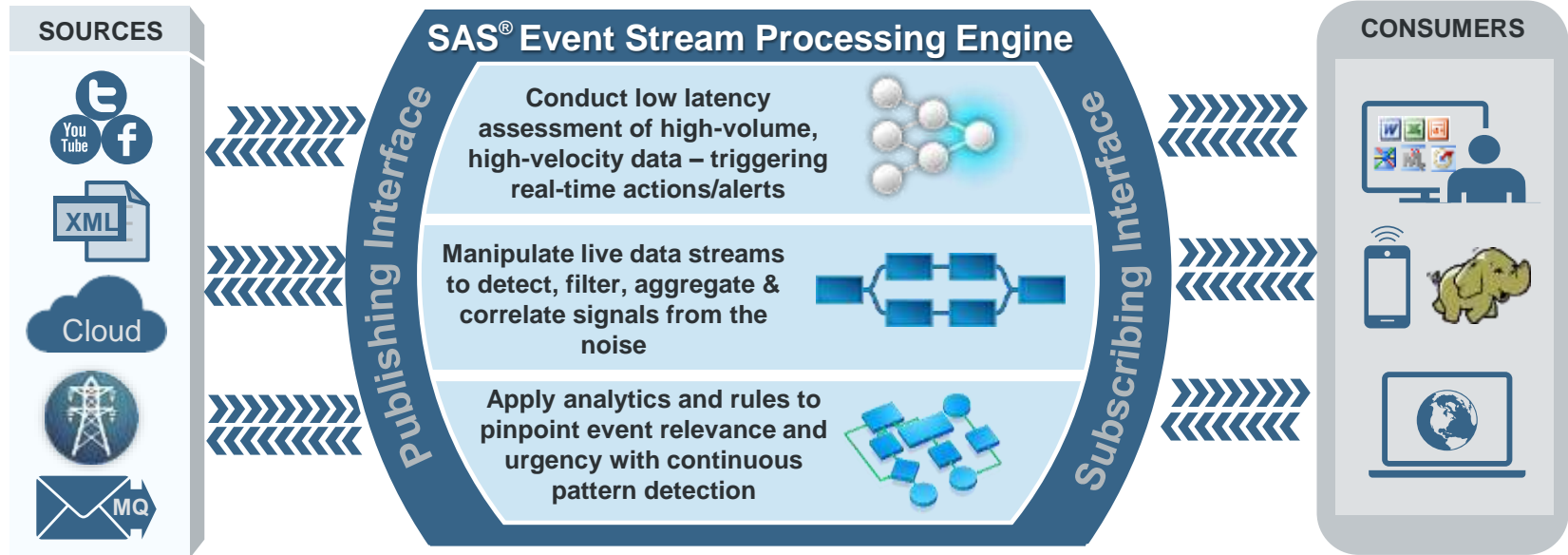
Pairwise Correlations for Table WORK.IRIS					
Column	Row	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
Sepal.Length	1	1.0000	-0.1176	0.8718	0.8179
Sepal.Width	2	-0.1176	1.0000	-0.4284	-0.3661
Petal.Length	3	0.8718	-0.4284	1.0000	0.9629
Petal.Width	4	0.8179	-0.3661	0.9629	1.0000

## SAS VISUAL ANALYTICS / STATISTICS



# PROCESAMIENTO COMPLEJO DE EVENTOS

# PROCESAMIENTO "ON-THE-FLY"





# PROCESAMIENTO COMPLEJO DE EVENTOS

## CASOS DE USO

### Optimización E-Commerce



- Análisis de uso de sitios Web
- Optimización de experiencia de usuario
- Publicidad y ofertas en tiempo real

### Detección de Fraudes



- Análisis de transacciones
- Alertas en tiempo real y manejo de casos.
- Correlación con comportamiento del cliente.

### Internet de las Cosas



- Sensores en tiempo real
- Detección de anomalías en tiempo real.
- Monitoreo de activos críticos
- Instrucciones en tiempo real

### Comunicaciones



- Alertas de mensajes en tiempo real.
- Ofertas y acciones en tiempo real.
- Diagnóstico y acciones

### Decision Management



- Decisiones operacionales centralizadas inmediatas.
- Directivas a empleados y sistemas en tiempo real.

### Mercados de Capital



- Cálculos y monitoreo continuo.
- Minimización del tiempo entre las operaciones y su reporte.



1. SAS da soporte al cliente en el diseño del cluster Hadoop y en establecer los criterios de éxito para las pruebas.
2. El cliente pone el cluster Linux (referencia 4 nodos, 8 cores, 128 GB RAM por nodo).
3. El cliente elige la distribución de Hadoop (Hortonworks / Cloudera).
4. SAS entrega el software por el período de prueba (referencia 3 meses):
  1. SAS Data Loader for Hadoop.
  2. SAS Visual Analytics / Statistics.
  3. SAS In-Memory Statistics for Hadoop.
  4. SAS Event Stream Processing.

5. Servicios profesionales de SAS instalan Hadoop y las soluciones SAS elegidas. Este es el único cargo al cliente a precio promocional.
6. Sugerencia de prueba:
  1. Formatos ORC/Parquet y SAS SPDE.
  2. Hive y SAS Data Loader for Hadoop.

**ESPERAMOS QUE  
NUESTRAS SOLUCIONES  
SEAN DE SU UTILIDAD**



**THE  
POWER  
TO KNOW.**