Conceptos Fundamentales de Arquitectura de Big Data

Computación de Alto Desempeño Diego Alberto Rincón Yáñez MCSc diego-rincon@javeriana.edu.co

Palabras de Koheleth Hijo de David, Rey en Jerusalem 200 D.C.

"No hay nada Nuevo bajo el Sol!"

Ecclesiastes Chapter 1 verso 9





Antes de empezar.....

¿Qué es Big Data?



¿De donde la moda del Big Data?

"Big Data is like teenage sex, everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it."

Dan Arieli,

Duke University Director

Center for Advanced Hindsight



¿Qué sabemos realmente de Big Data?

- Big Data es solo "data"
 - ¿Cual es la mejor forma de manejar nuestros datos?
- Big Data es una pieza en un rompecabezas.
 - ¿Cómo combinarla con la analítica existente para lograr un gran impacto?
- Big Data necesita unirse a la operación del negocio.
 - ¿Cómo podemos usar big data para crear mejores productos y servicios?

NO al botón de reinicio, no estamos iniciando desde cero.



¿Qué sabemos realmente de Big Data?

Nivel Organización

- Convertir la cultura de la organización en "data-driven".
- Usar preguntas de negocio, no de tecnología para guiar la arquitectura.
- Incrementar la autogestión de datos, para fomentar la decisión a partir de datos (Empoderar al Usuario).
- Habilitar un ecosistema que permita a los equipos "nadar" en los datos para descubrir.
- Hacer los datos y los resultados tan accesibles a la gente como sea posible.



¿Qué sabemos realmente de Big Data?

Nivel Tecnología

- Construir un modelo estructurado del negocio.
- Descubrir información acerca de clientes y productos.
 - ERP + CRM +Sistemas Legados
 - Logs web.
 - Correos + Redes Sociales
 - Voz + Texto
 - BI (Business Intelligence)
- Asegurar la información como un activo, controlando el acceso a los datos.
 - Proteger privacidad
 - Cumplir con regulaciones



Agenda

- Conceptos Básicos
- Fases del Proceso
- Arquitecturas de Big Data
- Tecnologías de Big Data
- Definición de Proyecto
- Casos de Estudio





Big Data & Data Science



Balanceo de Términos

- Analítica (Des, Diag, Pred y Pres)
- Data Science vs Big Data
- Arquitectura (Sw, IT, BD, Empr)
- B.D. (Architech, Engineer, Scientist)
- CDO (Chief Data Officer)
- (Machine o Deep) Learning
- Business Intelligence
- KPI (Key Performance Indicators)
- Open (Data vs Access)

- Metadata
- ETL vs ELT
- OLTP vs OLAP
- Data (Warehouse, Marts)
- Paralelo vs Distribuido
- Stream vs batch
- Clúster (HA vs HPC vs HTC)
- SQL y NoSQL
- Hadoop (HDFS+YARN+MR)



Big Data no es solo Hadoop

- Ejemplos donde Hadoop no es completamente aplicable
 - Cyber Seguridad, mercados de bolsa, información de sensores y tendencias en redes sociales (monitoreo en tiempo real).
 - Que pasaría si, su empresa tuviera muchos centros de datos, donde el ancho de banda es un problema.
 - Control de los datos/Seguridad (en sectores específicos).
- Más detalles en el módulo más adelante.



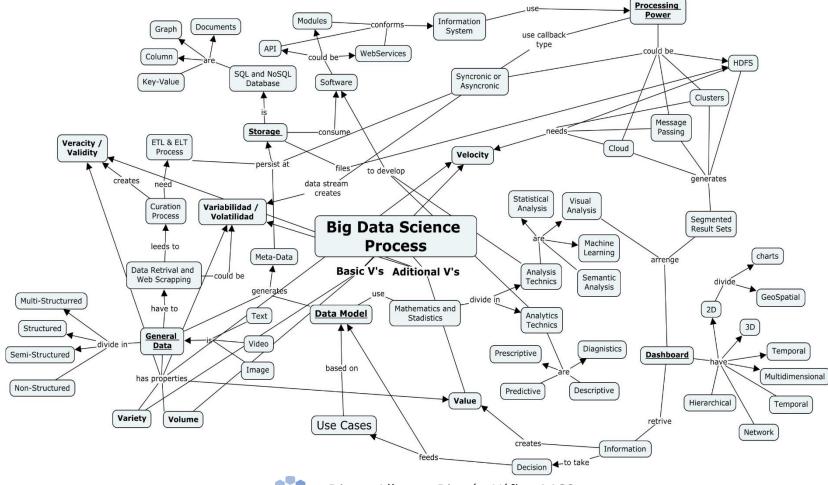
THE SIX

Phases Of A Project

- i. Enthusiasm
- 2. Disillusionment
- 3. PANIC
- A SEARCH FOR THE
- S. Punishment of the
- G. Praise & Honors for the honparheipants

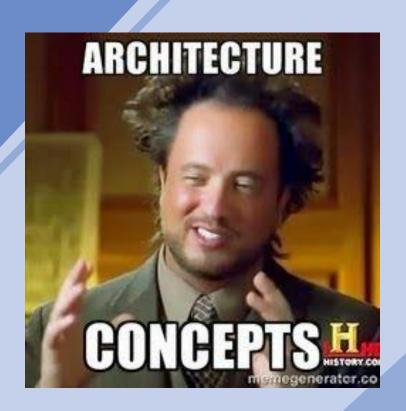
Fases del proceso





Diego Alberto Rincón Yáñez MCSc.

Twitter: @d1egoprog.



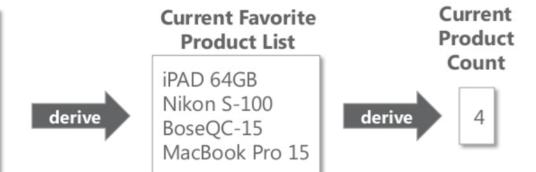
Conceptos de Arquitecturas de Big Data



Hoy Big Data es Procesamiento en Batch

Favorite Product List Changes

1.2.13	Add	iPAD 64GB
10.3.13	Add	Sony RX-100
113.13	Add	Canon GX-10
11.3.13	Remove	Sony RX-100
12.3.13 12.3.13 14.4.13	Add Add	Nikon S-100 BoseQC-15
15.4.13	Add	MacBook Pro 15
20.4.13	Remove	Canon GX10



Information => derived

Raw information => data



Diego Alberto Rincón Yáñez MCSc.

Twitter: @d1egoprog.

Operacional

Características

- Producida día a día
- Alto Volumen
- Baja Latencia
- Operaciones CRUD pequeñas (OLTP)
- Ejemplos
 - Datos de Cliente
 - Inventario
 - Compras



Características

- Múltiples Dominios (Negocios)
- Decisiones de Negocio
- Ajustado a la imaginación del analista.
- Optimizado para Mining, Ad hoc, batch, etc (OLTP)
- Ejemplos
 - Segmentación de Clientes
 - Pronósticos
 - Cambios de hábitos

Analítica



Arquitecturas en Números

	Operacional	Analítica
Latencia	1ms a 100ms	1min a 100min
Concurrencia	1000 a 100000	1 a 10
Tipos de Acceso	Lect / Escr	Lectura
Consultas	Selectivas	No Selectivas
Alcance	Operaciones	Histórico
Tecnología	NoSQL	MPP y MapReduce
Usuario Final	Cliente	Científico de datos



Principios de Arquitecturas Orientadas a Datos

- 1. Abrir la mente a los datos, todos, no solo los grandes, es critico extraer valor de todos.
- 2. Toda la tecnología es relevante, Data warehouse, descubrimiento y HDFS, son complementarias y buscar integración con el resto de la arquitectura empresarial, procesamiento de eventos, visualización y analítica.
- 3. Reemplazar la pila por una "cadena de suministro", un sistema que de respuesta a los datos no puede ser una pila, sino ajustarse como una cadena.



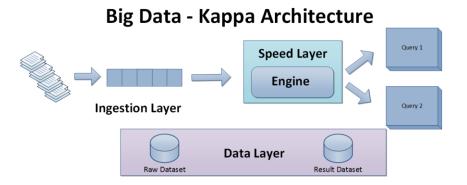
Principios de Arquitecturas Orientadas a Datos

- 4. Los consumidores de datos no son creados igual, los usuarios tiene diferentes habilidades, permisos y formas de recibir la información. Entregar información de la forma que la consumen.
- 5. Conserve lo que funciona, Establecer políticas alrededor de gobierno, calidad, seguridad y accesibilidad de los datos.
- 6. Empiece pequeño, Diseñar proyectos pequeños que den resultados rápido, no gastar tiempo en áreas de poco valor.

Arquitecturas estándares para Big Data



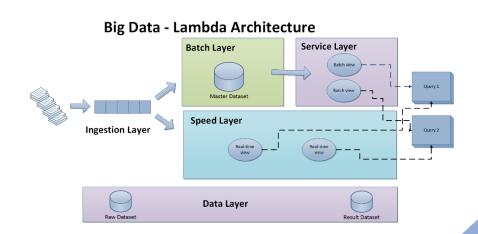
Arquitectura Kappa



- la capa de datos
- la capa de velocidad

Arquitectura Lambda

- la capa de datos
- la capa batch
- la capa de servicio
- la capa de velocidad





Taxonomías de Arquitecturas de Big Data

¿Quiénes usan arquitecturas de referencia de big data?



Taxonomías de Arquitecturas de Big Data







¿Quiénes usan arquitecturas de referencia de big data?

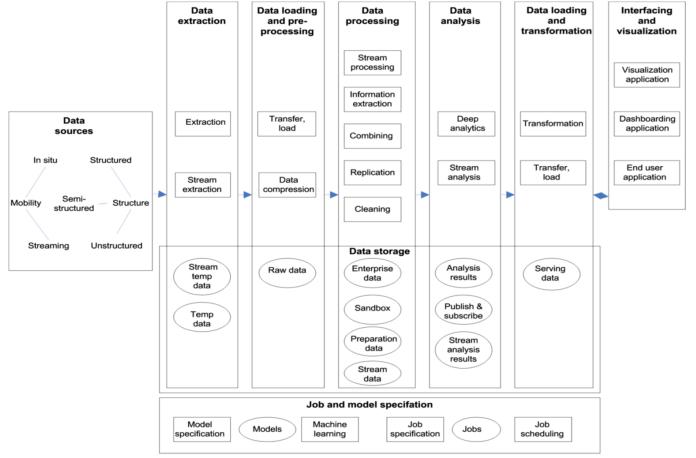








Taxonomía de las Arq. de referencia de Big Data





Diego Alberto Rincón Yáñez MCSc.

Twitter: @d1egoprog.

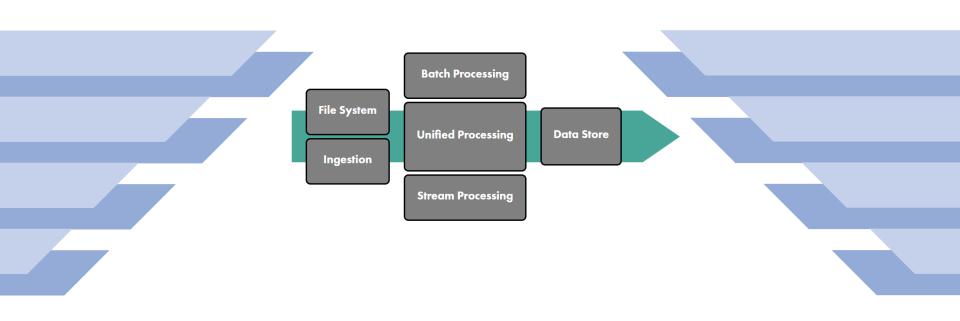


Tecnologías de Big Data



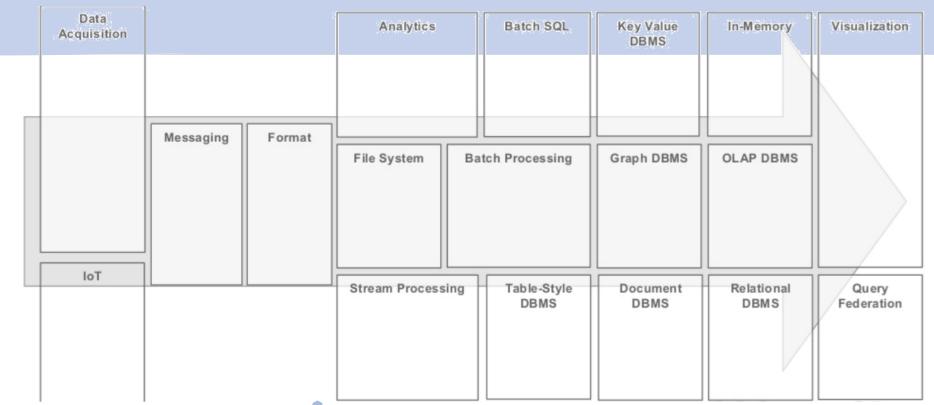
BIG DATA LANDSCAPE 2017







Plantilla de Tecnologías



Diego Alberto Rincón Yáñez MCSc.

Twitter: @d1egoprog.

Exploradas en este módulo

- Processing
 - YARN
 - MapReduce
- Analytics
 - SPARK
- Database
 - MongoDB



Definición de proyecto



Descripción del Problema

Utilizando una herramienta de sniffering, provista en ambientes de tipo UNIX (pueda correr en segundo plano de forma parametrizable). El equipo de investigación capturará (Streaming) el tráfico de la red de la Universidad Javeriana, con el fin de realizar un análisis del mismo.

Para lo anterior, deberán utilizar las arquitecturas, herramientas, conceptos vistos en clase y en la maestría para, construir una pequeña plataforma enfocada en el análisis tráfico de red en tiempo real y posible predicción del comportamiento de la red.



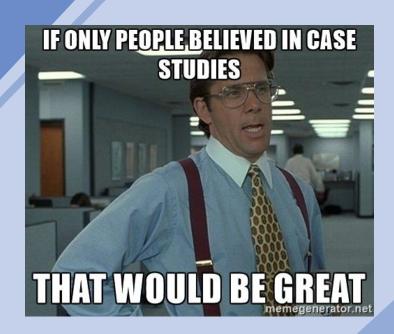
Reglas de Entrega

- Usar máquinas asignadas.
- Sustentación final 22 de Noviembre
- Nota por todo el equipo (Con aportes personales).
- Mínimos de Entrega
 - Debe tener instalado y hacer uso de:
 - ✓ Spark
 - ✓ MongoDB en HA.
 - ✓ HDFS + YARN + MapReduce.
 - Documentación del proyecto (Formatos en Blackboard)
 - Documentación de(los) datasets.
 - Documentación de Arquitectura.
 - Flujo de los datos y transformaciones en Información.

Calendario de Entregas

- Entrega Parcial 1 (Jueves 1 de Noviembre)
 - Presentación (30 minutos)
 - Alcance de Solución
 - Descripción de alcance analítico (# Indicadores).
 - Definición inicial de Arquitectura
 - Blueprint de tecnologías
- Entrega Final (Jueves 22 de Noviembre)
 - Documentación Completa
 - Prototipo Funcional + Montaje de Servicios
 - Presentación de resultados (10 Minutos)





Casos de Estudio



Definición de Equipos

- FB
- NF
- TW
- IN



Reglas de Presentación

- Introducción de la Problemática
- Contexto de la Solución
 - Funcionalidades
 - Alcance
- Explicación de la Arquitectura
 - Restricciones
 - Ventajas
- Flujo de datos y transformaciones.
- Blueprint de Tecnologías
 - Ecosistema e integración.
 - Descripción y funcionalidad (por Tecnología)
 - Entradas/Salidas (por Tecnología)
- Aprendizajes y Factores a Resaltar
- Conclusiones

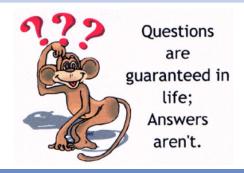
1 de Noviembre (20 Minutos)





Diego Alberto Rincón Yáñez MCSc.

Twitter: @d1egoprog.



¿Preguntas?



Diego Alberto Rincón Yáñez MCSc. Twitter: @d1egoprog.