

# Programa AWS Big Data – Analytics Specialist



“Big Data, el nuevo petróleo”

## Reflexión de la noche

« Un esfuerzo más y lo que iba a ser un fracaso se convierte en un éxito; no existe el fracaso si nos esforzamos cada vez mas y mas »

**ANITA QUEVEDO()**



**Data Engineer Sr;**



;

**Directora / Docente de  
Tecnologías Cloud**



;

**Docente de  
Big Data -  
Spark**



Cloud & Data Analytics

**Fundadora  
Comunidad**

Estudiante apasionada, profesora y profesional de TI con 7 años de experiencia en proyectos de BI, BA Big Data y Cloud Computing, y cursando el Master en Data Management en U. de Barcelona

**Trabajé en:**



an NTT DATA Company

# Metodología de enseñanza



- ★ **Complementar la teoría con casos reales**
- ★ **Discusiones tecnológicas**
- ★ **Laboratorios con ejemplos útiles y reusables**
- ★ **Material complementario para profundizar**
- ★ **Comunicación 360**

# Módulos del programa



- ★ Overview herramientas
- ★ Collect
- ★ Storage
- ★ Procesamiento
- ★ Analítica
- ★ Seguridad
- ★ Visualización

# Te queremos conocer...



Nos interesa saber tus expectativas, dudas y experiencia.

1. ¿Donde trabajas cuáles son tus principales funciones?
2. ¿Qué conoces de Big Data?
3. ¿Que esperas del programa?

Cosas a tener  
en cuenta



Se creará un grupo  
de wsp para estar  
todos comunicados



El material estará  
disponible en  
**Google Drive**

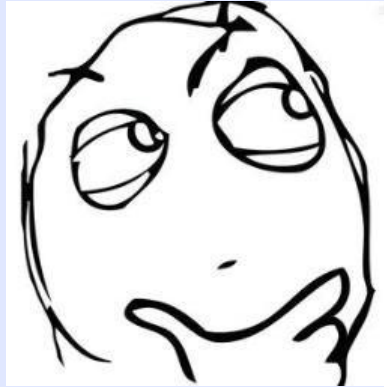
# **Sesión 1: Introducción**



## Sesión 1: Objetivos

- ✓ Definir Big Data
- ✓ Identificar algunas fuentes de Big Data
- ✓ Listar ejemplos de uso de Casos de Big Data
- ✓ Describir el ecosistema de Big Data
- ✓ Review del big data “pipeline” y las herramientas disponibles para cada fase.

# ¿Qué es el Big Data?



# ¿Qué es el Big Data?

La colección y análisis de grandes volúmenes de data:



Responder preguntas



Crear ventaja competitiva

# Cuando los dataset son grande , resultan difíciles para

- Almacenar
- Organizar
- Analizar
- Mover
- Compartir

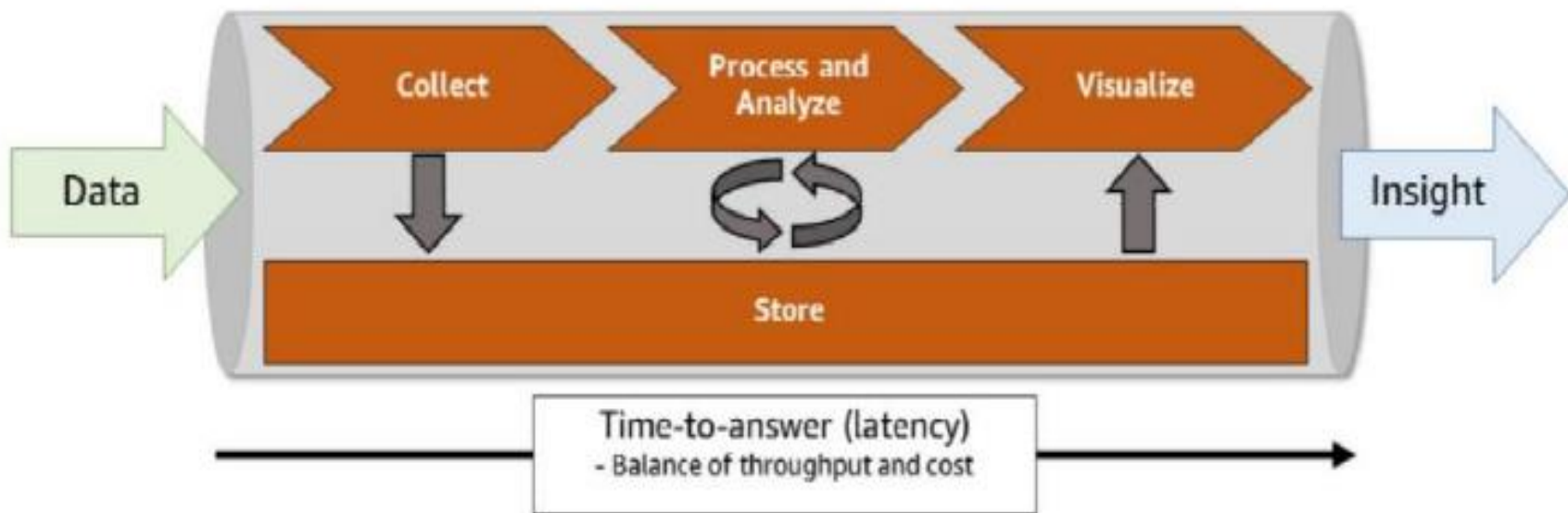


Big Data puede generarse a partir de muchas fuentes diferentes, que incluyen:

- ❑ Redes sociales
- ❑ Servicios bancarios y financieros.
- ❑ Servicios centrados en la web.
- ❑ Búsquedas científicas y de documentos
- ❑ Registros médicos
- ❑ Registros web



# **El Big Data “Pipeline”**



## Collect

**Near Real-time**  
Amazon Kinesis Firehose

**Data Import**  
Amazon Import/Export Snowball

**Message Queuing**  
Amazon SQS

**Web/app Servers**  
Amazon EC2

## Store

**Object Storage**  
Amazon S3  
Amazon Glacier

**Near Real-time**  
Amazon Kinesis Streams

**RDBMS**  
Amazon RDS

**NoSQL**  
Amazon DynamoDB

**Search**  
Amazon CloudSearch

**Internet of Things (IoT)**  
Amazon IoT

**On Demand Mode**  
DynamoDB

**Double Peak Performance  
with IO Volumes**  
Amazon EBS

**Transactional API**  
DynamoDB

## Process and Analyze

**Hadoop Ecosystem**  
Amazon EMR

**Near Real-time**  
AWS Lambda  
Amazon Kinesis Analytics

**Data Warehousing**  
Amazon Redshift

**Machine Learning**  
Amazon SageMaker

**Elastic Search Analytics**  
Amazon Elasticsearch Service

**Process and Move Data**  
AWS Data Pipeline  
AWS Glue

**Embed Interactive Dashboards**  
QuickSight

**Analytics for Java Apps**  
Amazon Kinesis

**Ad Hoc Analytics**  
Amazon Athena

## Visualize

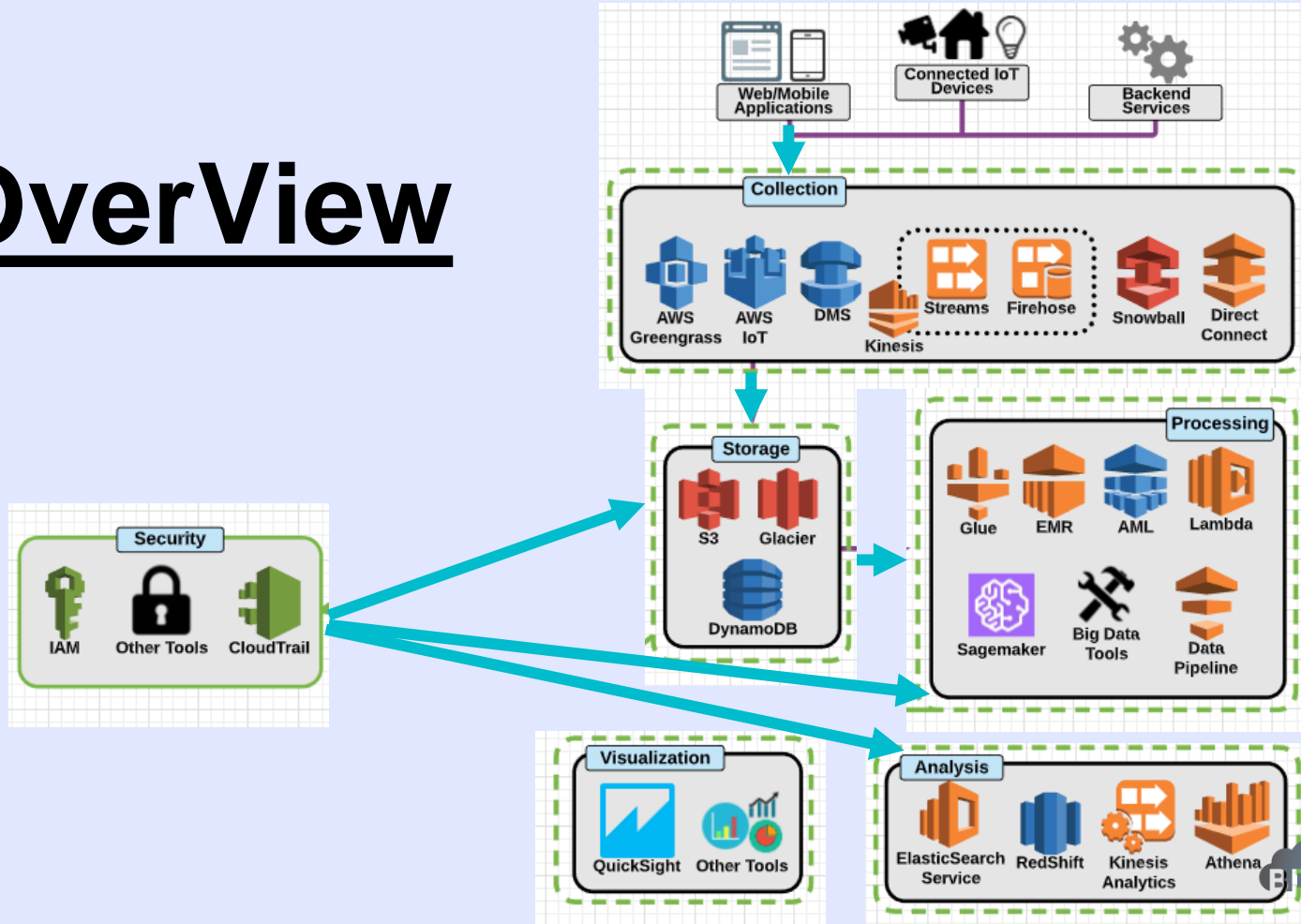
**Business Intelligence and Data  
Visualization**  
Amazon QuickSight

**Elastic Search Analytics**  
Amazon Elasticsearch Service

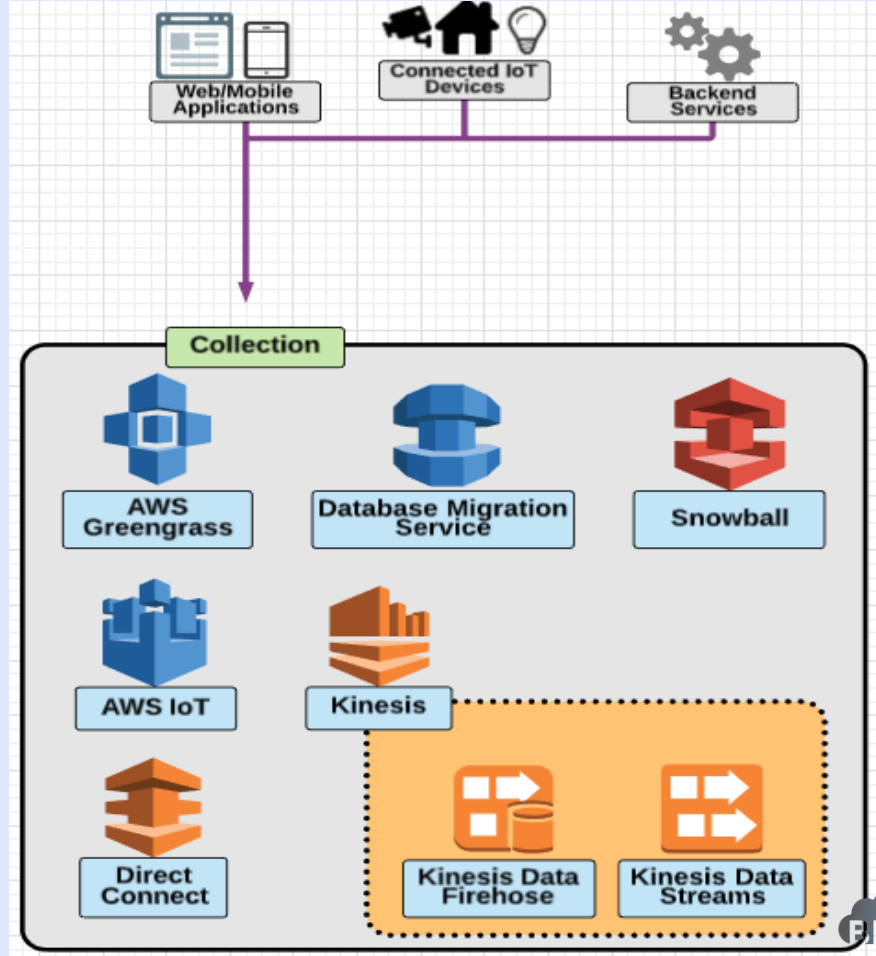
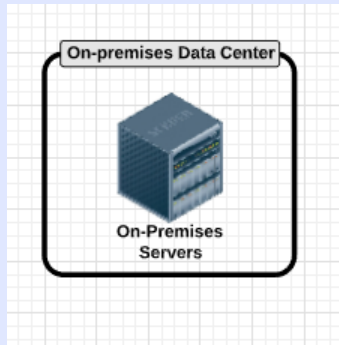


# “OverView

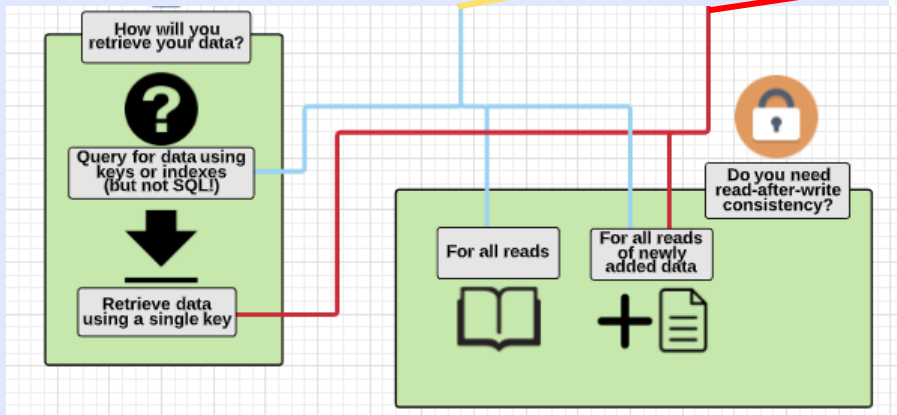
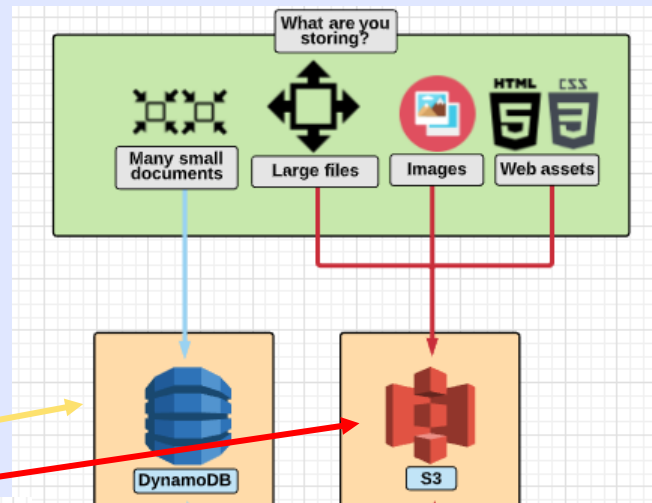
”



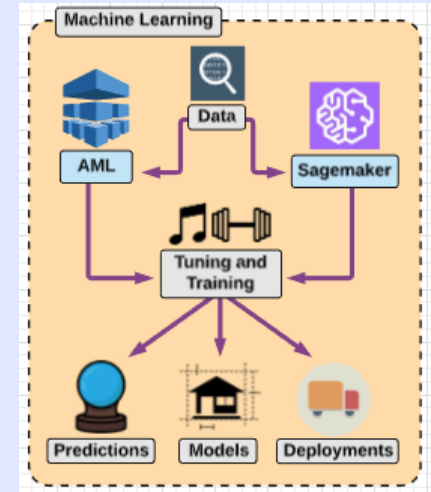
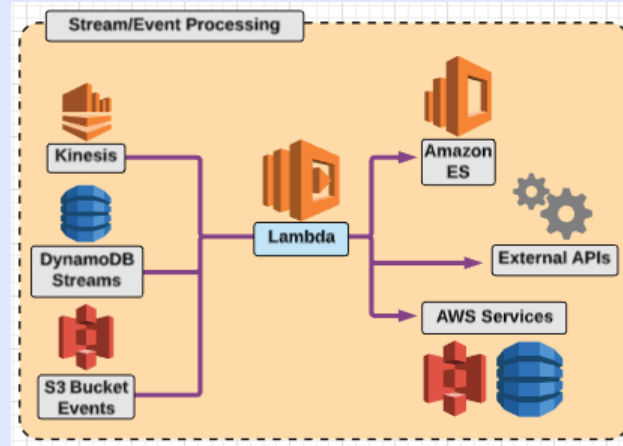
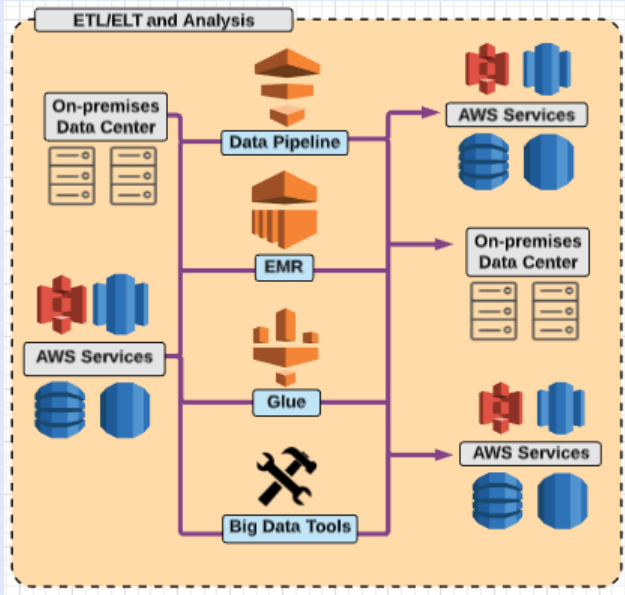
# Colección



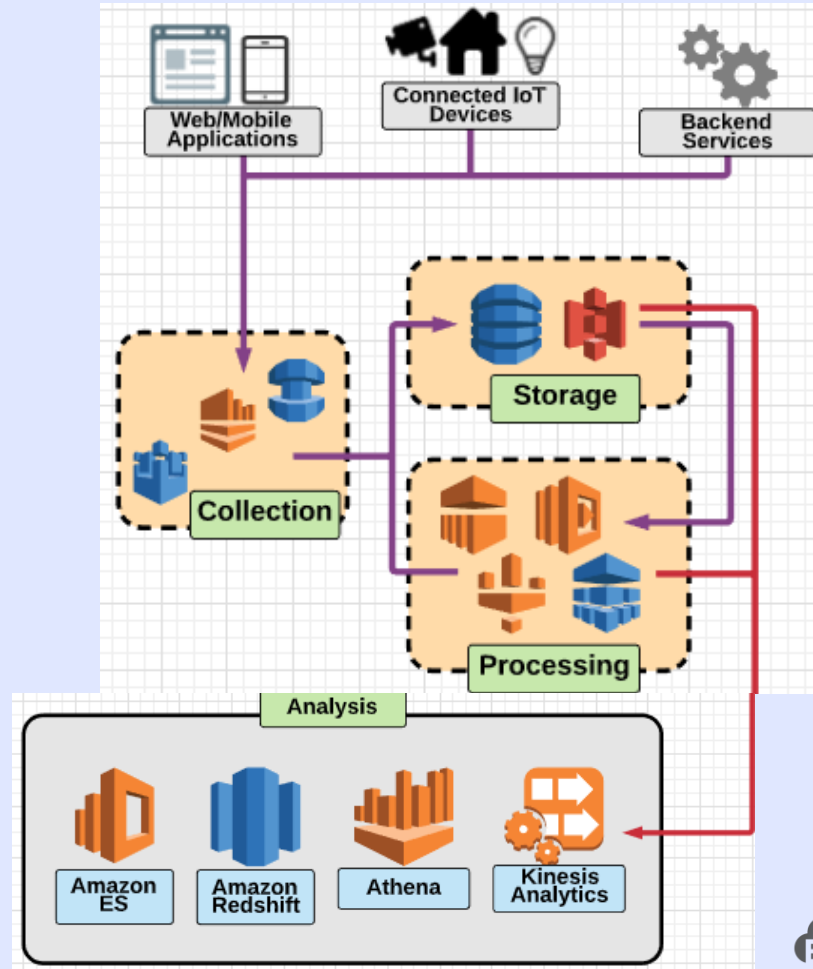
# Storage



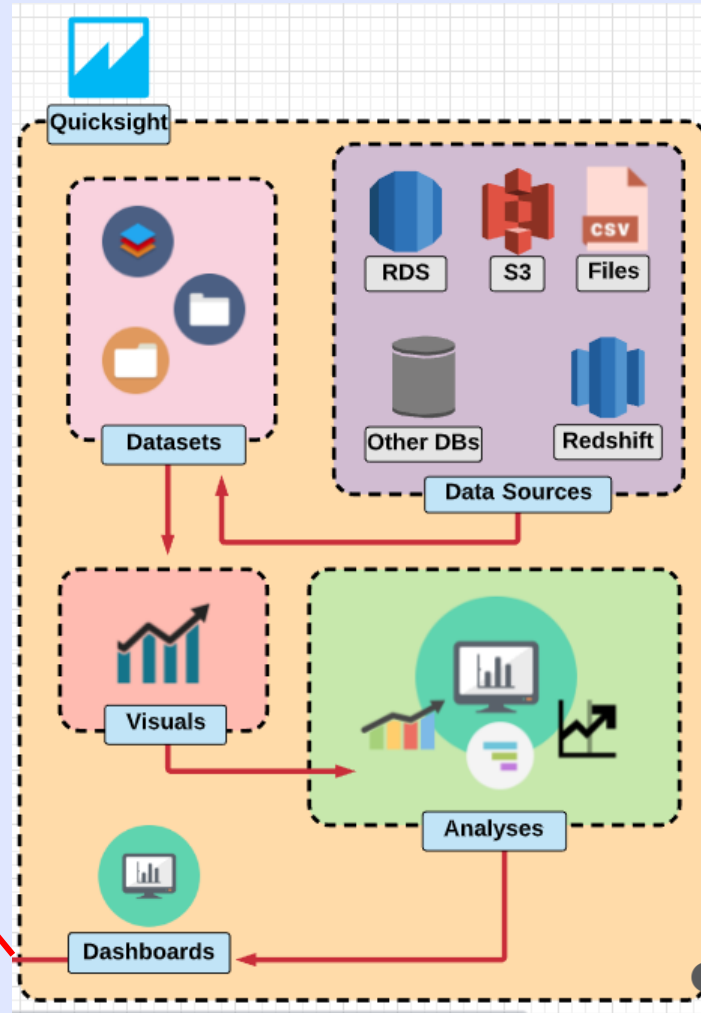
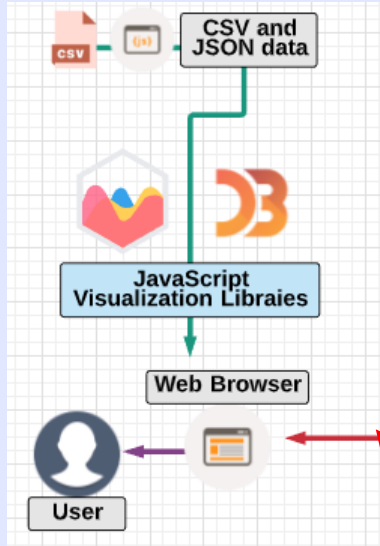
# Procesamiento



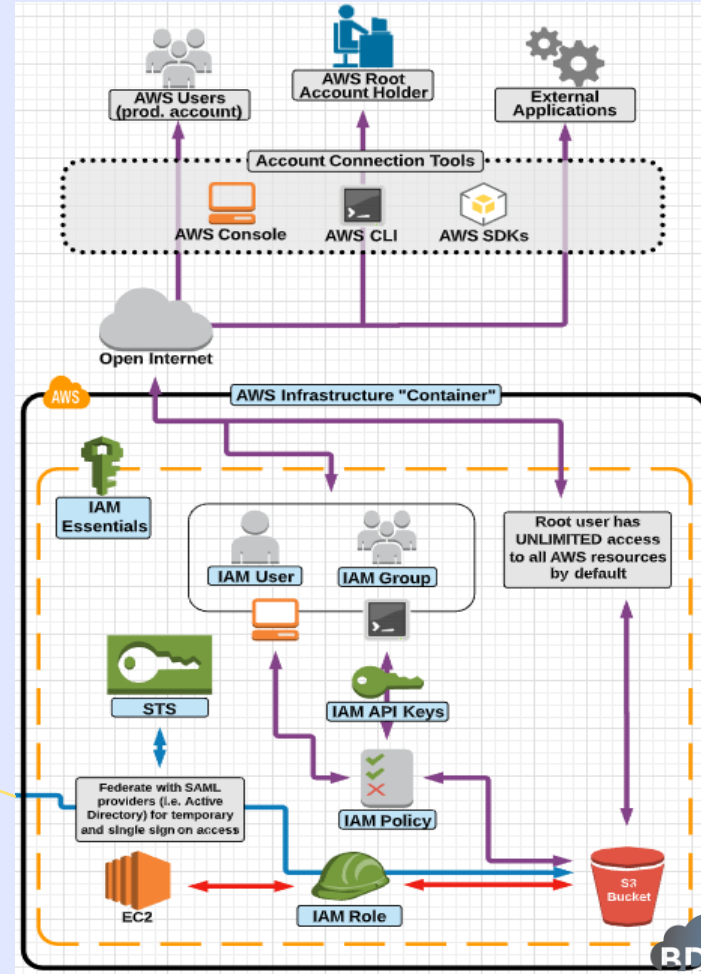
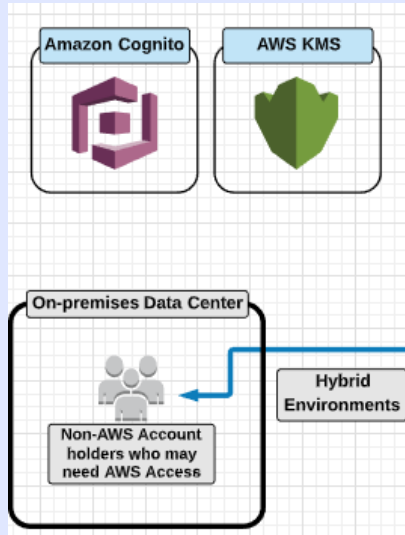
# Análisis



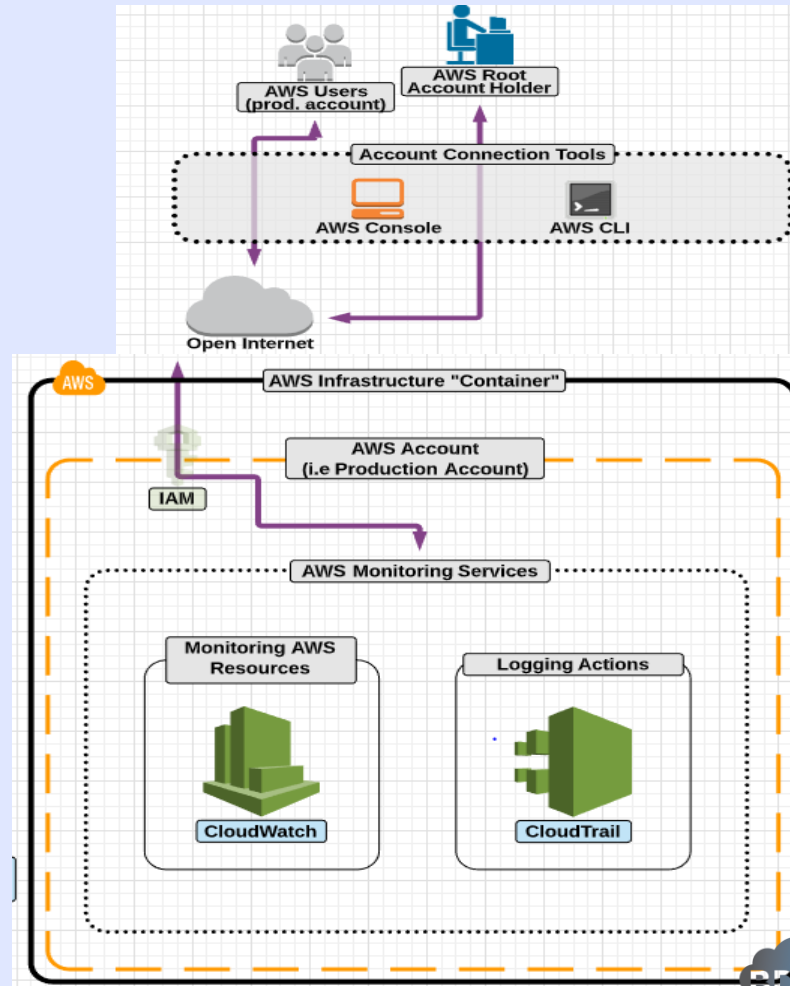
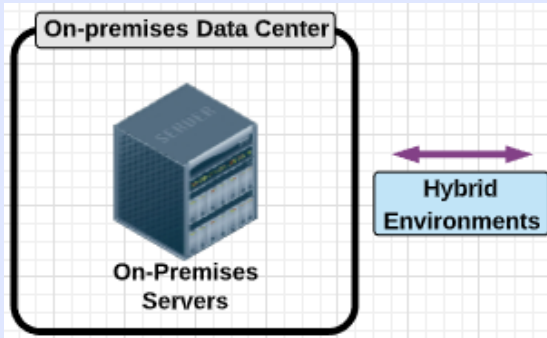
# Visualización



# Seguridad

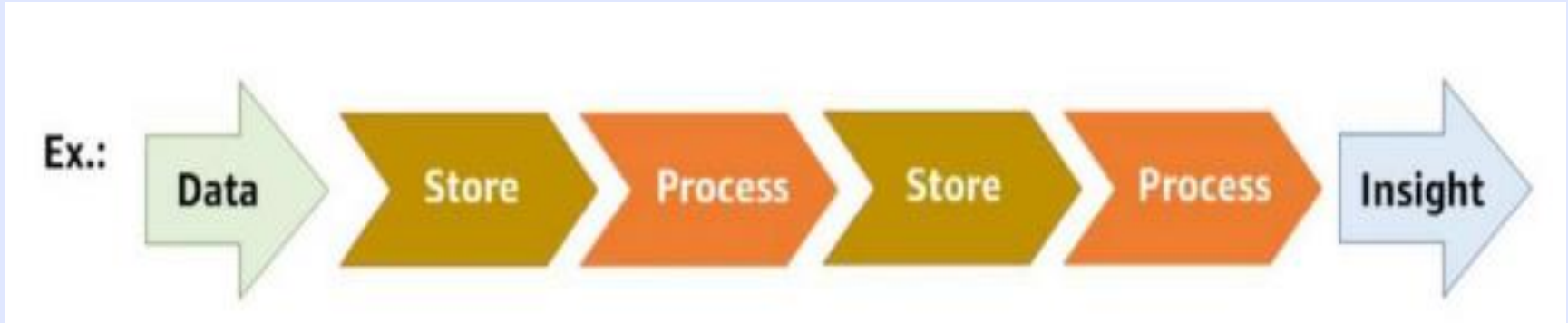


# Monitoreo





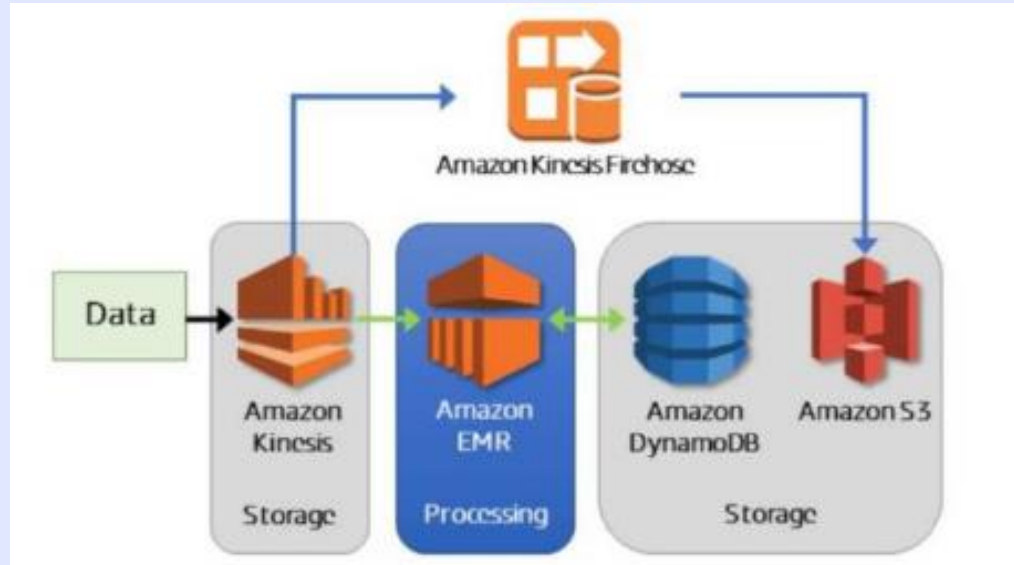
## Desacoplar su bus de datos



Haga coincidir la herramienta con la tarea requerida.  
esto da como resultado una arquitectura más tolerante a fallas y  
generalmente resulta en un tiempo mejorado para responder

## desacoplar su bus de datos: Flexibilidad

con un bus de datos desacoplado, puede hacer que múltiples aplicaciones de procesamiento lean o escriban en múltiples almacenes de datos



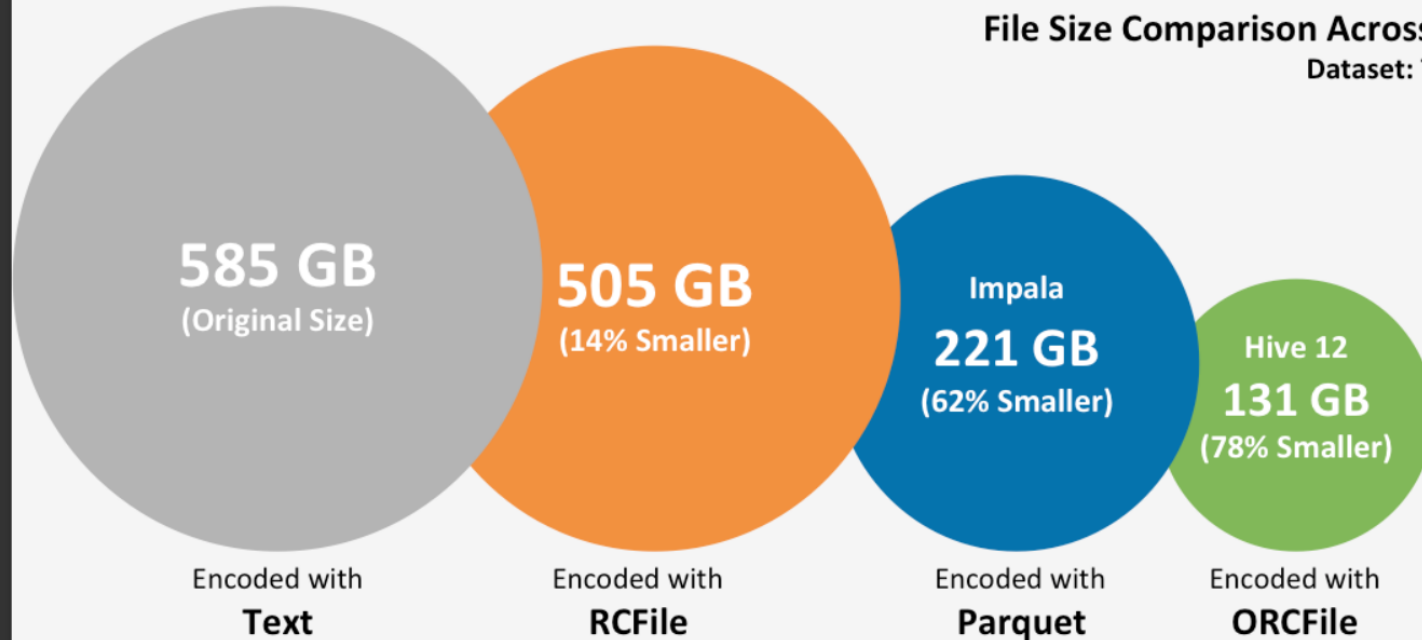
# Data Lake options on AWS

---

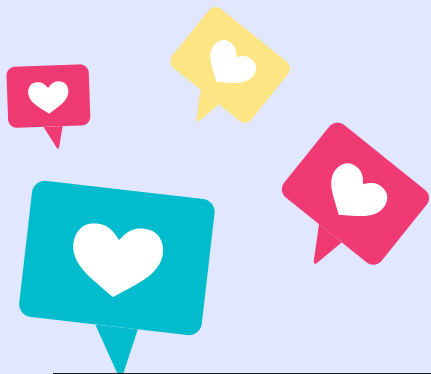
Storage	Processing	AWS-Managed Solution	Vendor-Managed
HDFS	Spark	AWS EMR (HDFS+Spark)	EC2 + Vendor Solution
S3	Spark	AWS EMR (Spark Only)	EC2 + Vendor Solution
S3	Serverless	AWS Athena	Serverless + Vendor Solution

## File Size Comparison Across Encoding Methods

Dataset: TPC-DS Scale 500 Dataset



- Larger Block Sizes
- Columnar format arranges columns adjacent within the file for compression & fast access



## ¿Qué tecnología de procesamiento de datos debo usar?

	Amazon Redshift	Impala	Presto	Spark	Hive
Latencia de consultas	Baja	Baja	Baja	Baja	Media (Tez) – Alta (MapReduce)
Durabilidad	Alta	Alta	Alta	Alta	Alta
Volumen de datos	1,6 PB máx.	Nodos aprox.	Nodos aprox.	Nodos aprox.	Nodos aprox.
Administrada	Sí	Sí (EMR)	Sí (EMR)	Sí (EMR)	Sí (EMR)
Almacenamiento	Nativo	HDFS/S3A*	HDFS/S3	HDFS/S3	HDFS/S3
Compatibilidad con SQL	Alta	Media	Alta	Baja (SparkSQL)	Media (HQL)

Baja Baja Baja Media Alta  
Latencia de consultas (Baja es mejor)



MUCHA TEORIA  
A PRACTICAR



# Gracias

