

# Clase 4 IA - BigData

## Tipos de datos

- Datos Estructurados: A menudo numerosos etiquetas almacenadas en un marco estructurado de columnas y filas.
- Datos Semiestructurados: Organizados libremente en categorías utilizando etiquetas meta.
- Datos No Estructurados: Información con mucho texto que no está organizada en un marco o modelo claramente definido.

## ETL vs ELT

### ETL

- Extract: Se obtienen los datos desde diversas fuentes (api, db).
- Transform: Los datos extraídos son procesados y transformados fuera del sistema destino. Se limpian, formatean y estructuran de acuerdo a las reglas de negocio.
- Load: Los datos transformados se cargan en un almacén.

### ELT

- Extract: Los datos son extraídos de las fuentes, igual que en ETL
- Load: Los datos sin transformar se cargan directamente en el sistema destino. > Es un concepto alimentado por la nube
- Transform: La transformación ocurre dentro del sistema de destino, aprovechando su capacidad de procesamiento.

## Categorías

### Bronce

Tener los datos en bruto, por ejemplo, cuando hacemos web-scraping. a veces es json, csv. mas o menos estructurados.

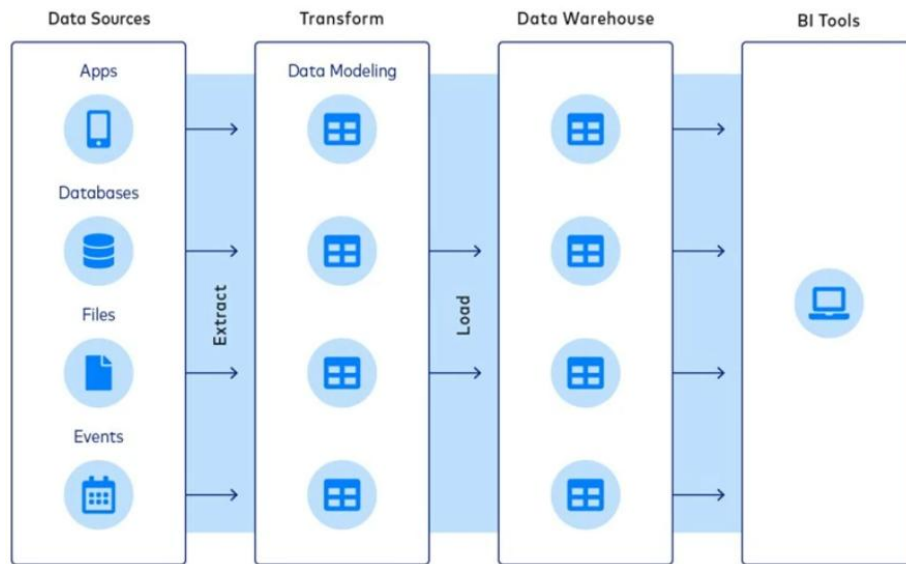


Figure 1: Imagen de referencia

## Plata

es el procesamiento de los datos de tipo bronce. simplemente una pequeña modificación de los datos que los deje listos para hacer una transformación mas completa.

## Oro

Es el dato más refinado. Esto es luego de hacer agregaciones (porque necesitan mucho computo). Este es el nivel desde donde debería consumir el sistema

Herramientas que se usan: > Azure, Apache spark.

## Webscrapping

Es un proceso mediante el cual extraemos datos de la web de forma medianamente automatizada.

## Nota

escribir la diferencia entre **Data Factory** contra **Data Bricks**.