

Curso de

# Data Warehousing y Modelado OLAP

Edison Yepes

# ¿Qué es BI y Data Warehousing?

# 2021 This is What Happens In An Internet Minute



Creado por: LoliLewis OfficiallyChadd

# ¿Qué es BI?

“Habilidad para **entender** la interconexión de los hechos presentes como **medio para descubrir guías** que nos permitan llegar a metas futuras, lo que significa mejorar nuestra capacidad de hacer juicios **predictivos** para el acierto en la **toma de decisiones.**”

[Hans Peter Luhn, 1958]

“BI usa datos de **ayer y hoy** para tomar **las mejores decisiones** acerca del **mañana**”

[Scheps, 2008]

“Habilidad para **transformar** los **datos** en información, y la **información** en **conocimiento**, de forma que se pueda optimizar el proceso de **toma de decisiones** en los negocios”

[Mora, 2007]

“Conjunto de **estrategias y herramientas** enfocadas a la administración y creación de **conocimiento** mediante el **análisis de datos** existentes en una organización o empresa.”

[Howard Dresner (Grupo Gartner), 1989]

# ¿Qué es Data Warehousing?

“

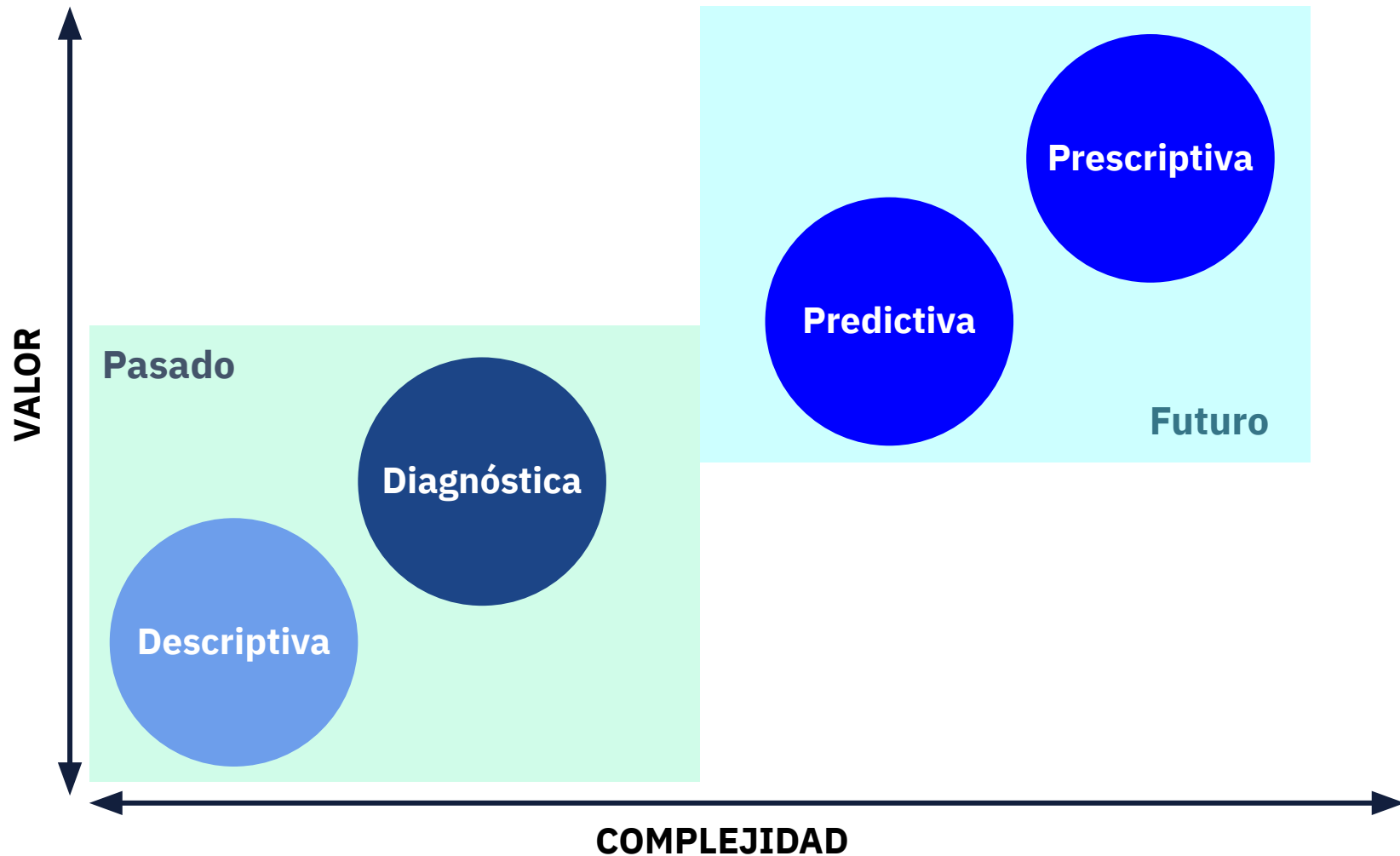
**Es un proceso, no un producto,  
para ensamblar y administrar  
datos de diversas fuentes con el  
fin de obtener una visión única  
y detallada de una parte  
o de toda una empresa.**

”

*Devlin, 2011*

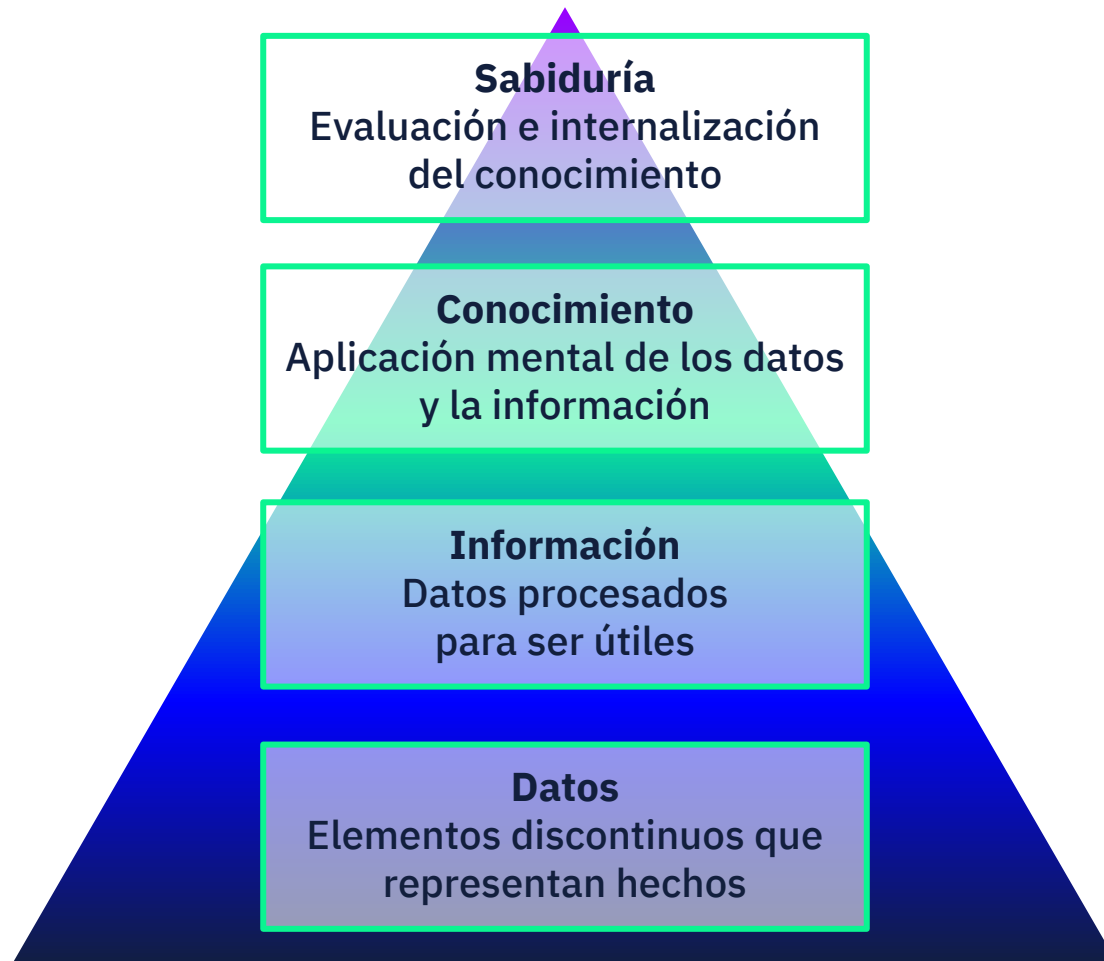
# Niveles de analítica y jerarquía del conocimiento

# Niveles de analítica

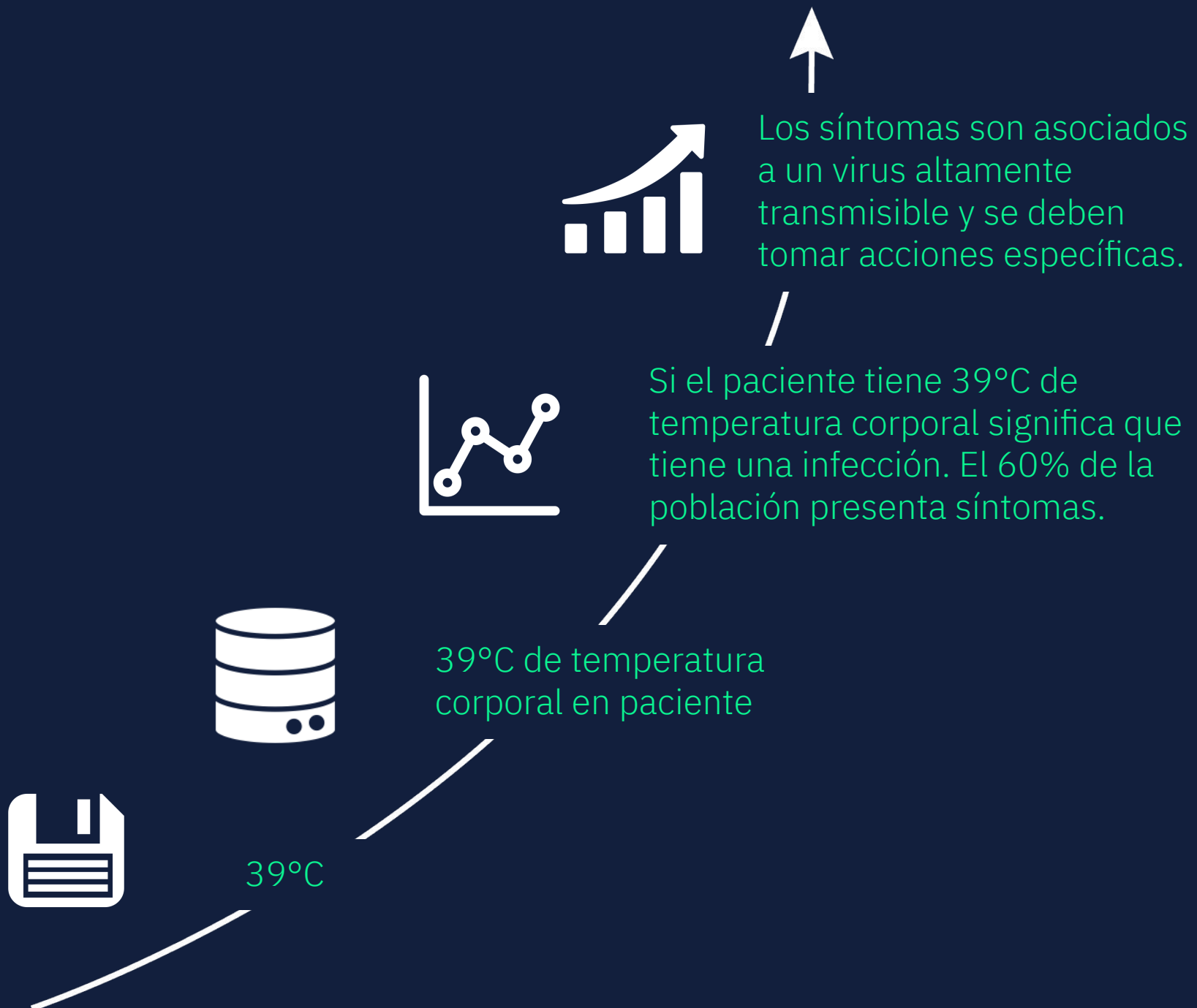




# Jerarquía del conocimiento



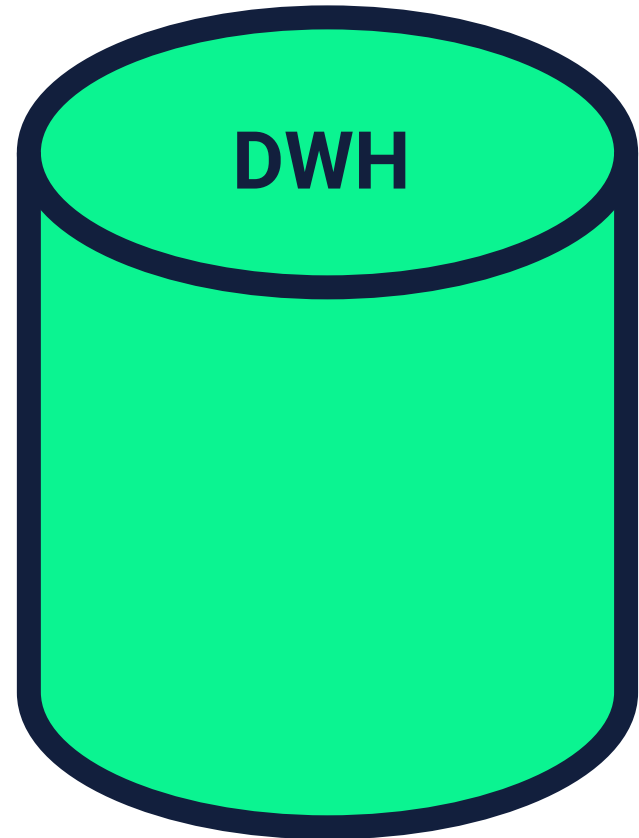
**Pirámide D-I-K-W**



# **Data Warehouse, Data Mart, Dimensiones y Hechos**

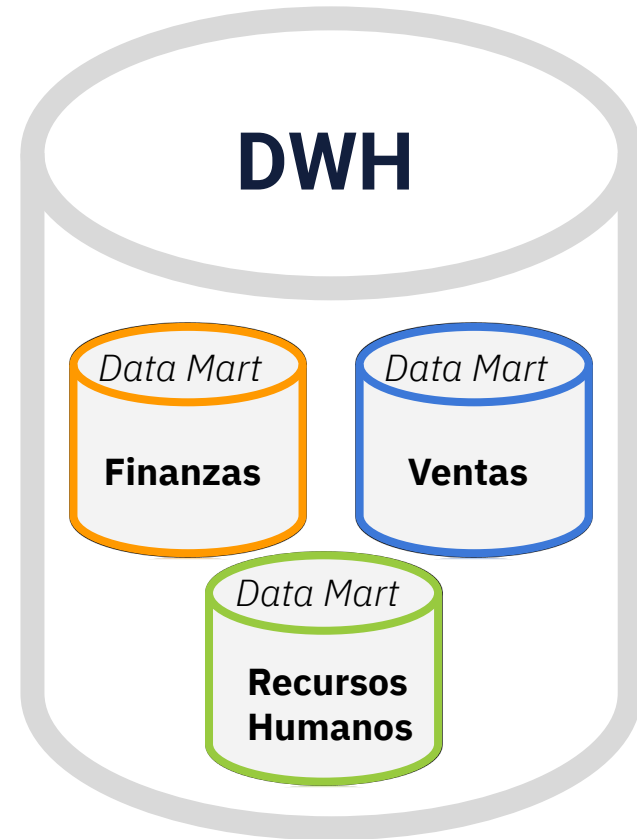
# Data Warehouse (DWH)

- Base de datos que contiene información de muchas fuentes diferentes.
- Los informes creados a partir de un Data Warehouse son usados para tomar decisiones.



# Data Mart

- Segmento del DWH orientado a un área específica del negocio.
- Contienen información sumariada para el análisis en una unidad de la organización.



# Dimensiones

- Describen los procesos del negocio.
- Diferentes actores en los procesos del negocio.
- Entidades del negocio.



# Hechos

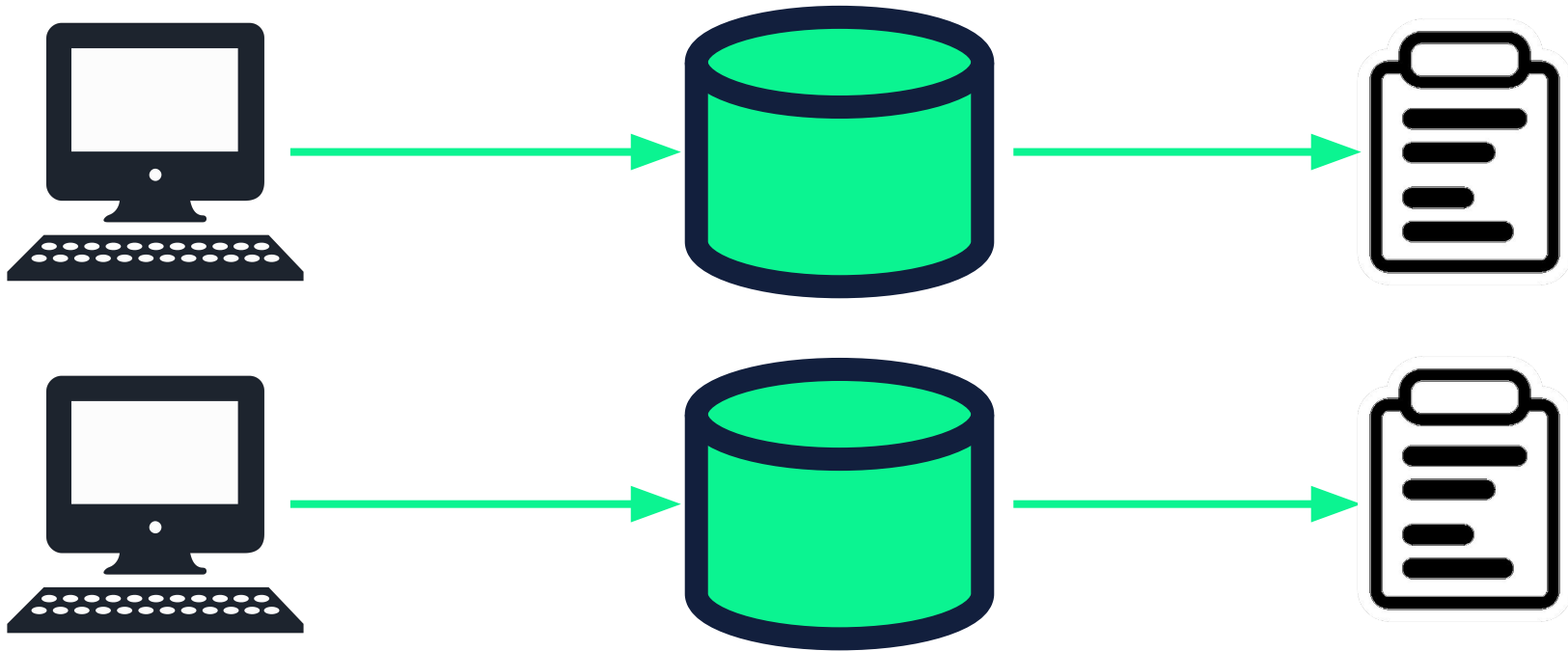
- Es la información cuantitativa de un proceso de negocio.
- Se denominan **Medidas** o **Métricas**.



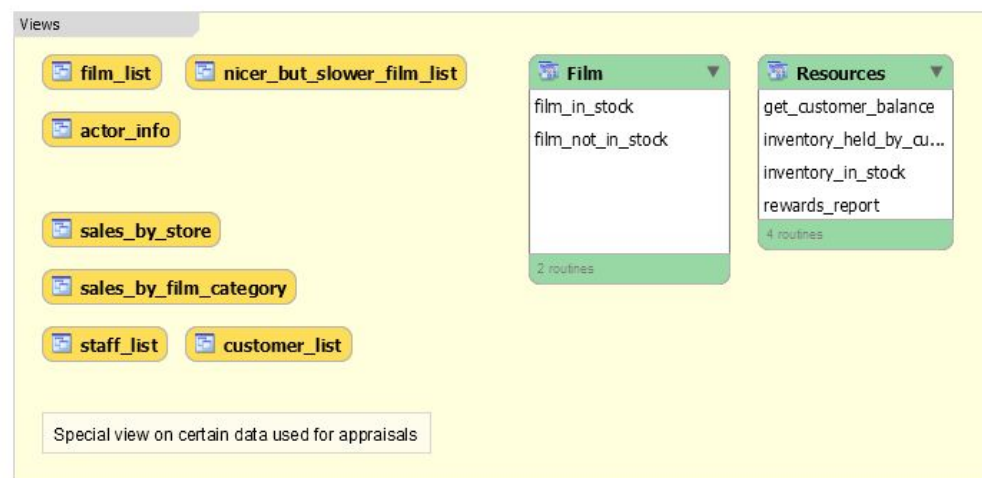
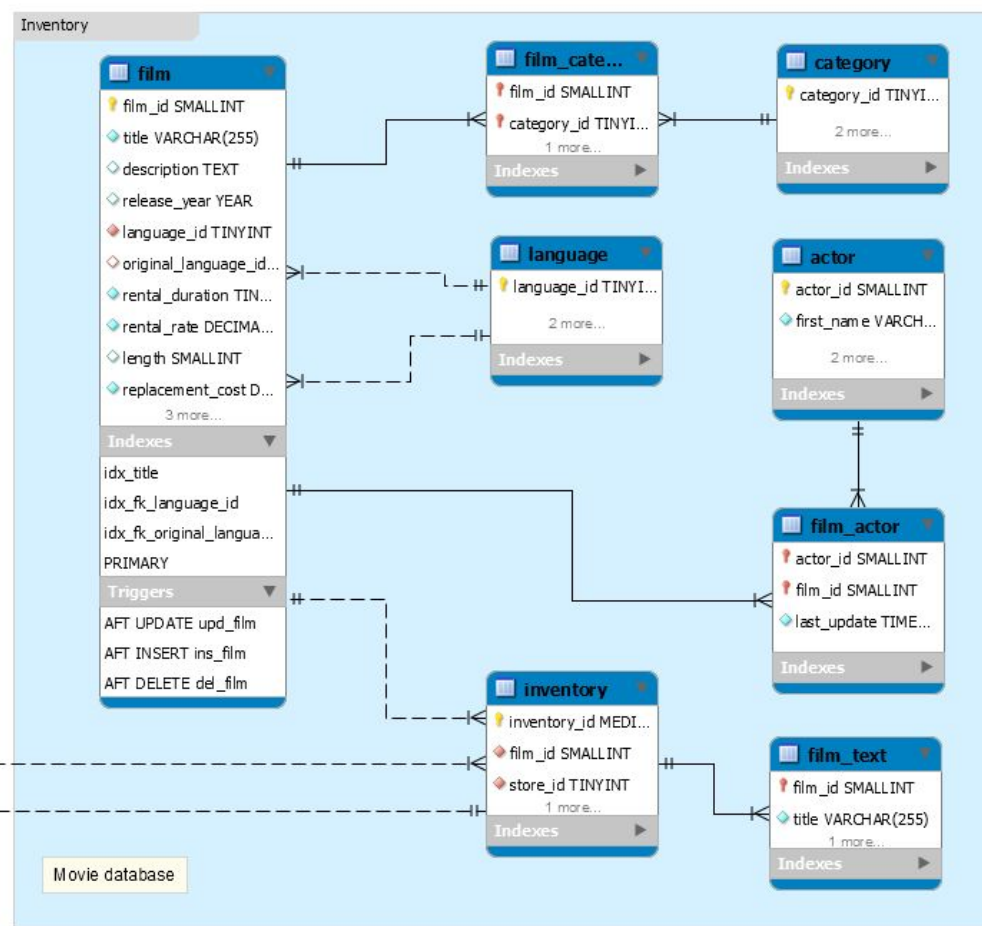
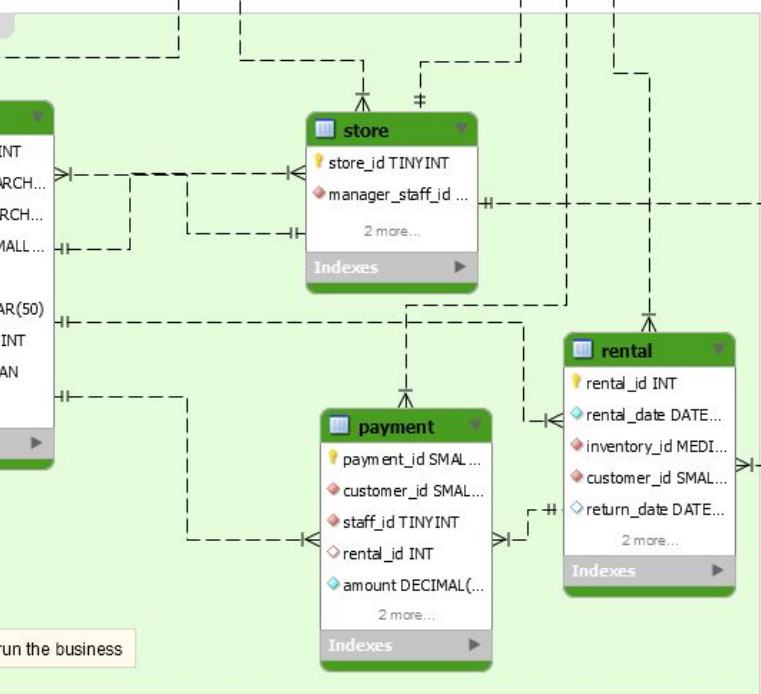
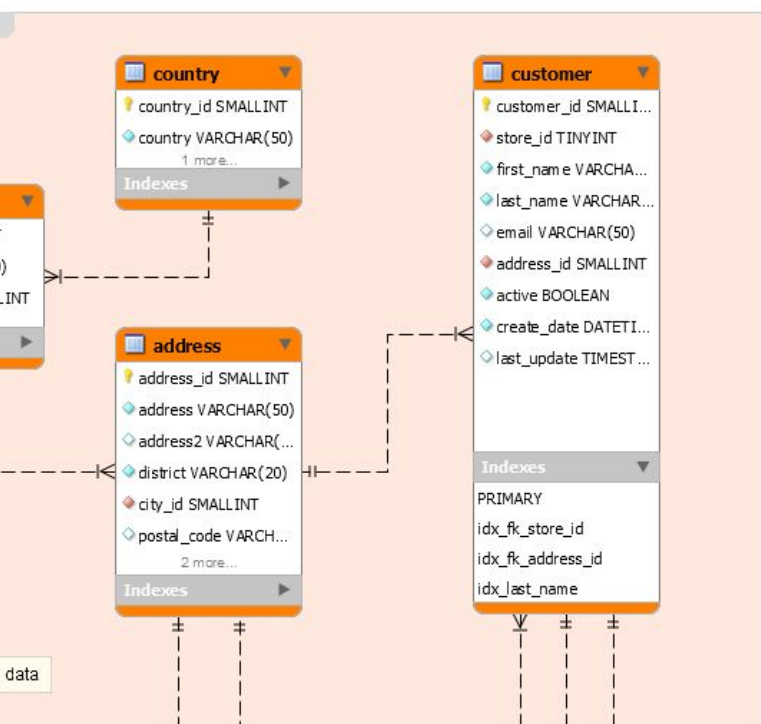
# OLTP vs. OLAP



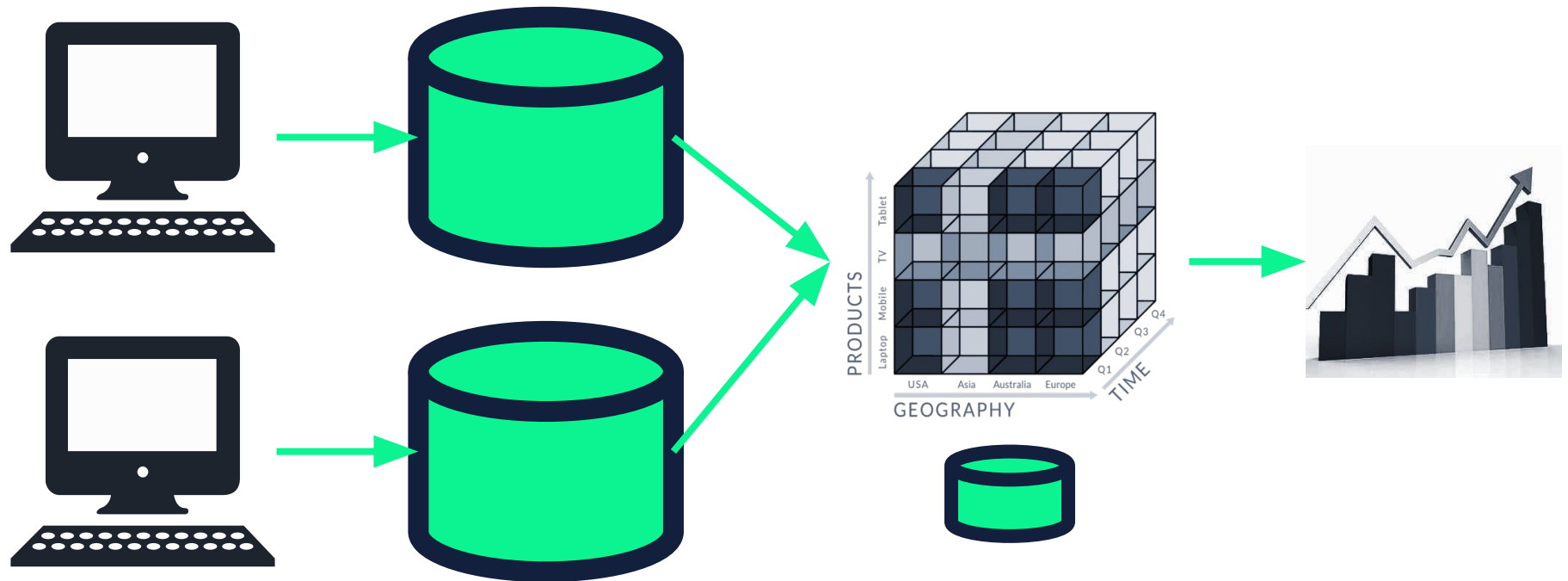
# OLTP



**Procesamiento de  
Transacciones En Línea**

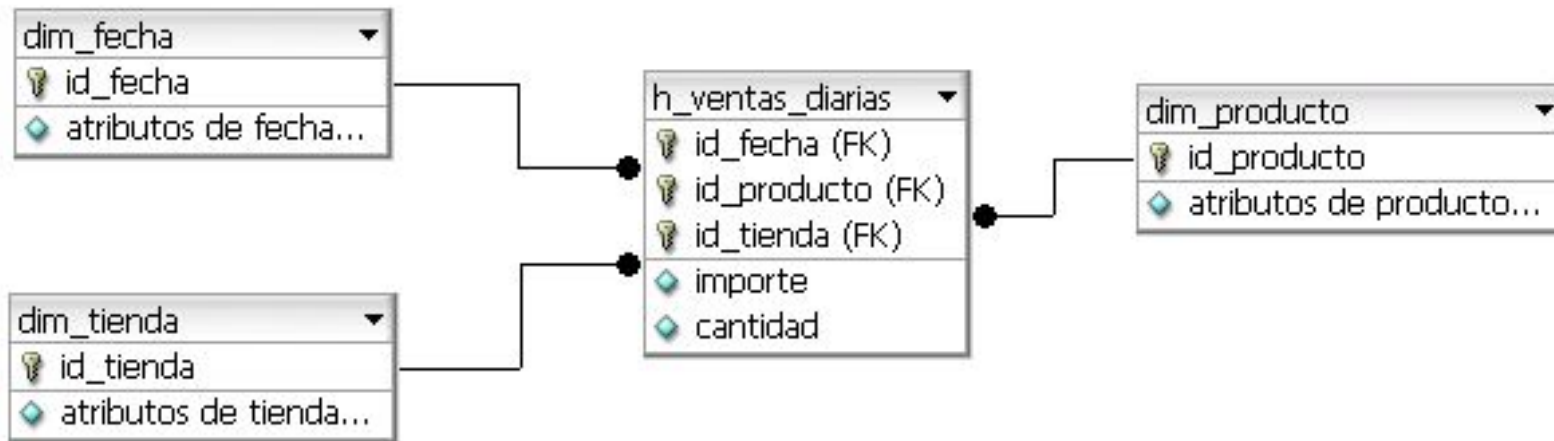


# OLAP



**Procesamiento Analítico En Línea**

# OLAP



# OLTP vs. OLAP

## Sistema Operacional (OLTP)

- ¿La factura ### fue cancelada?
- ¿Cliente que compró el producto X el día de hoy?

---

## Sistema de Bodega de Datos (DW)

- ¿Producto más vendido en el año 2022 por línea de producto?
- ¿Comparación de las ventas vs. el presupuesto, mes a mes, por tienda?

# OLTP



**ORACLE**

# OLAP



**Azure Synapse Analytics**



**Google BigQuery**

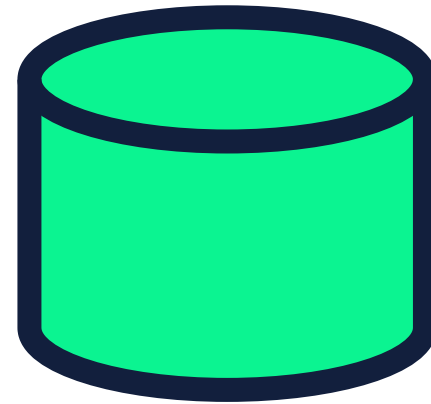
OLTP	OLAP
Diseñado para soportar las transacciones comerciales.	Diseñado para apoyar el proceso de toma de decisiones.
Data es volátil.	Data NO es volátil.
Data detallada.	Datos resumidos.
Modelado E-R.	Modelado dimensional.
Procesamiento de transacciones.	Procesamiento analítico.
Alta concurrencia.	Baja concurrencia.

# Metodologías de Data Warehouse

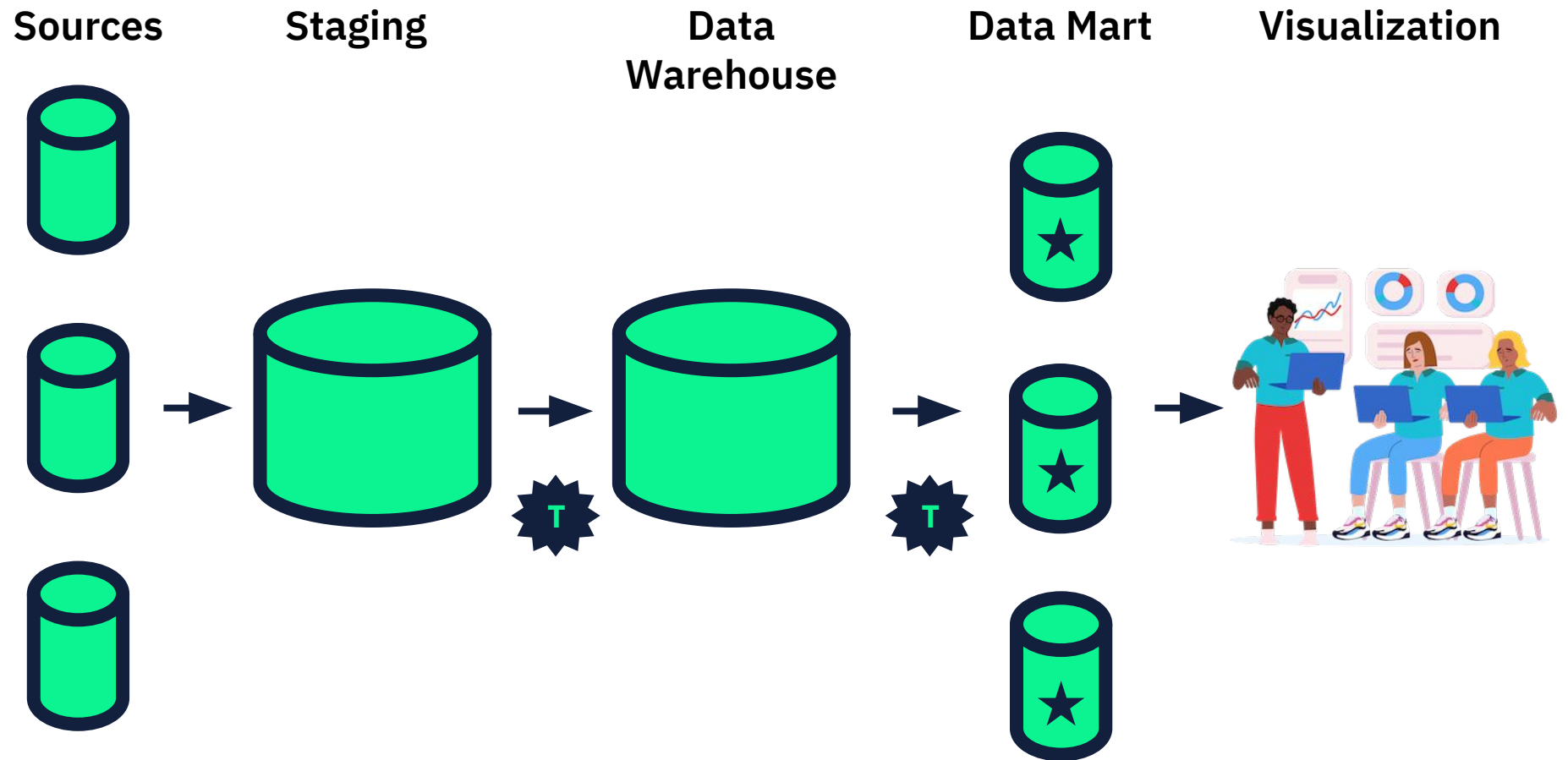


# Metodologías de DWH

- Sources
- Dimensional Models
- Visualization

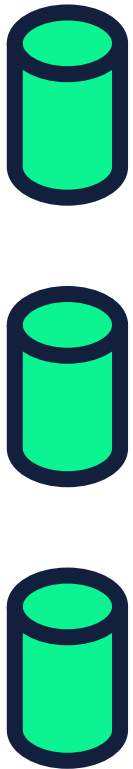


# Bill Inmon

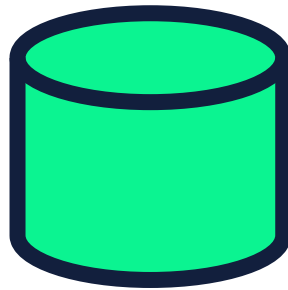


# Ralph Kimball

Sources



Staging



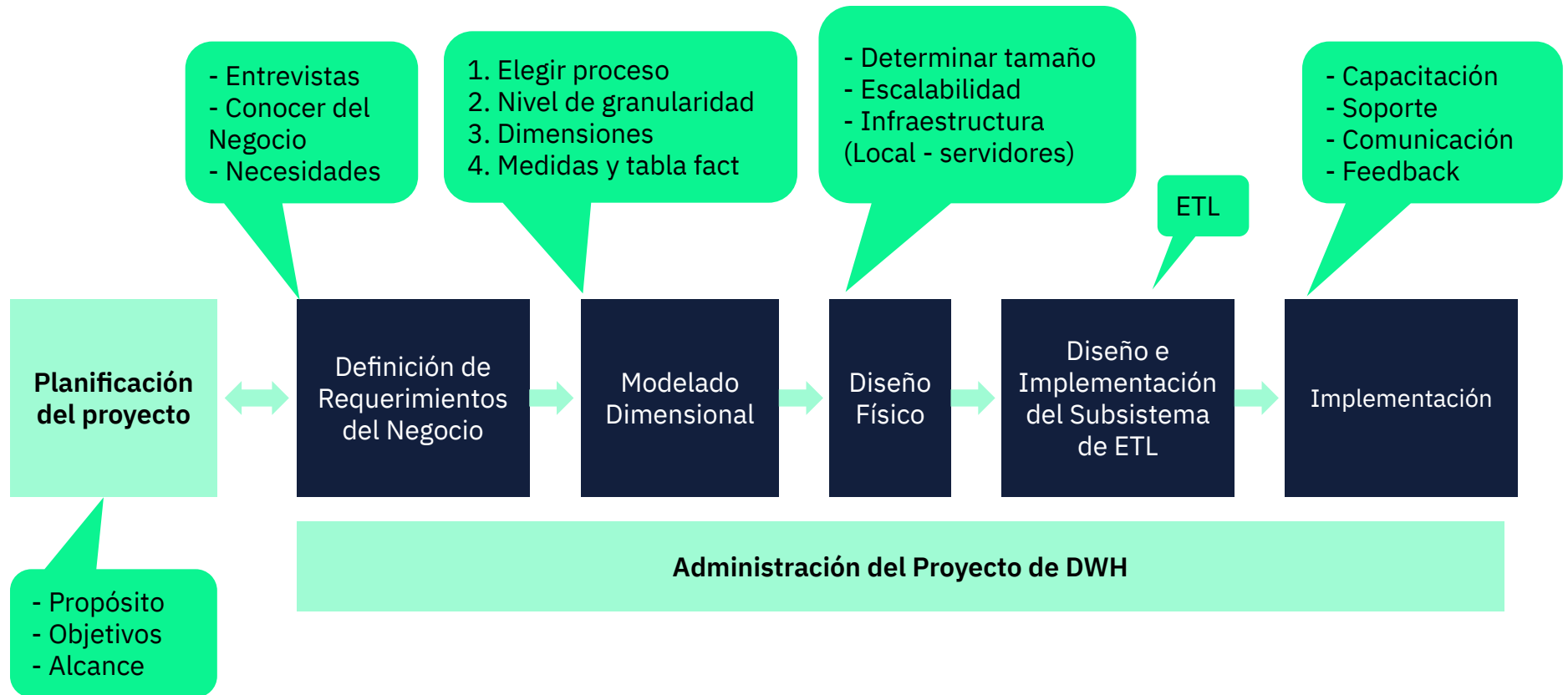
Data Mart



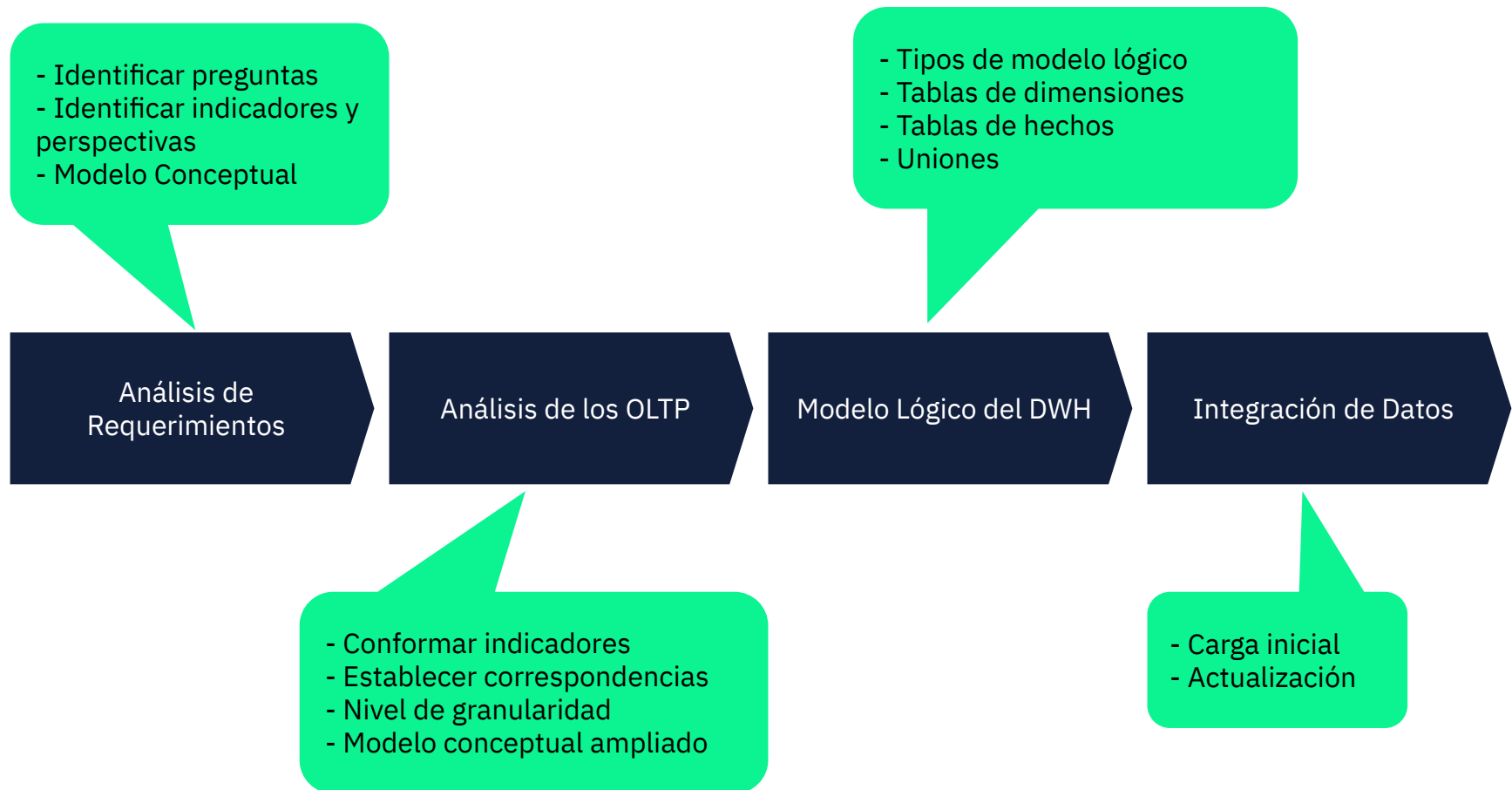
Visualization



# Ralph Kimball - fases

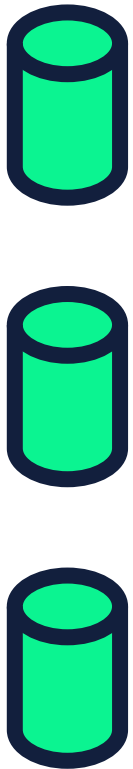


# Hefesto

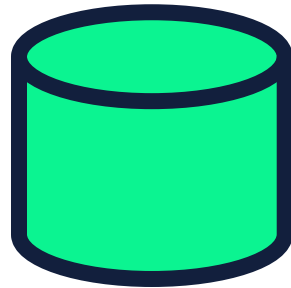


# Personalización

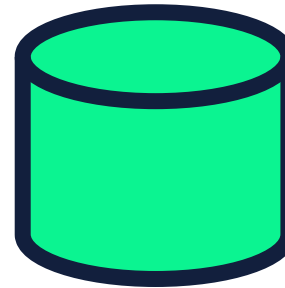
Sources



Staging



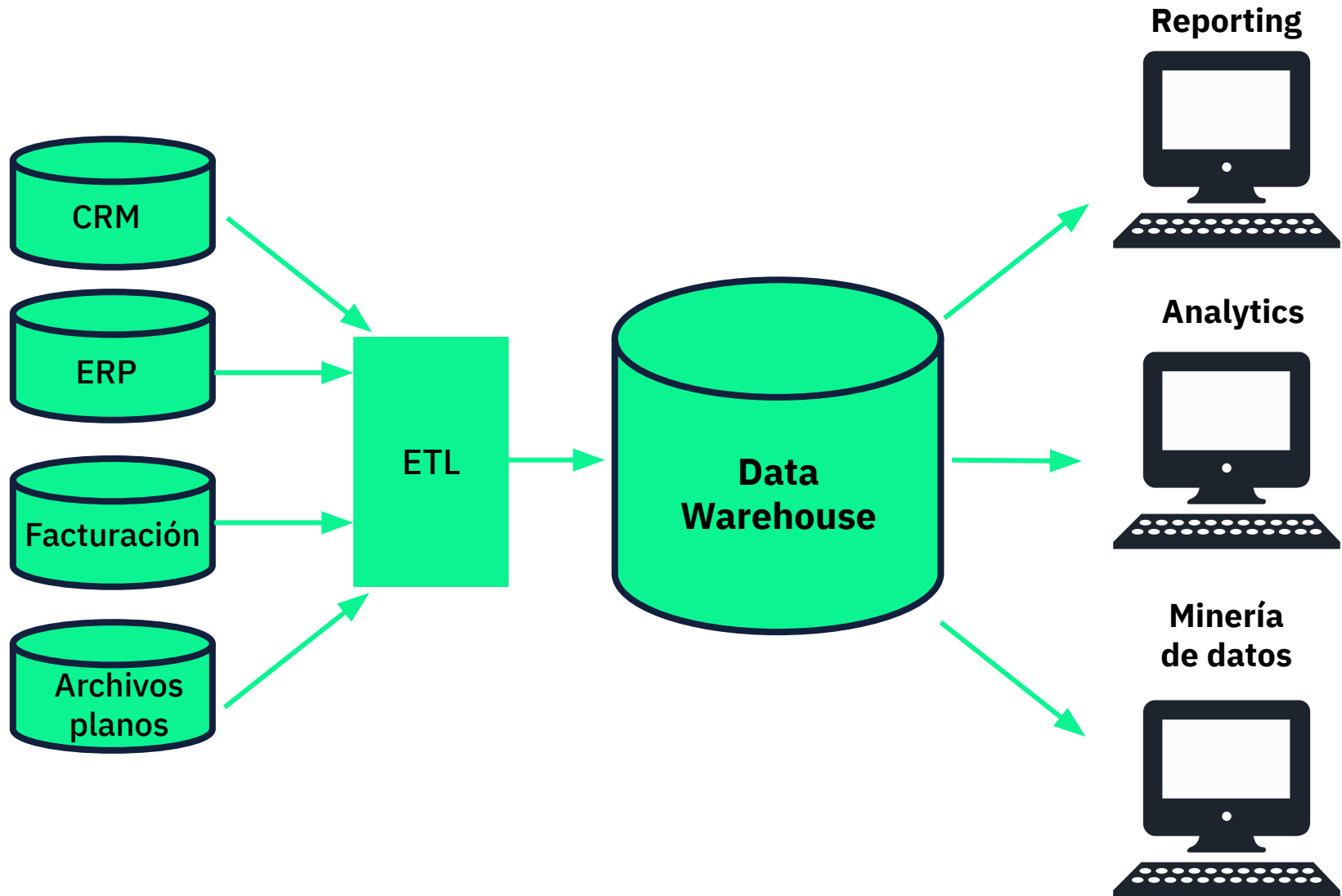
Data  
Warehouse



Visualization

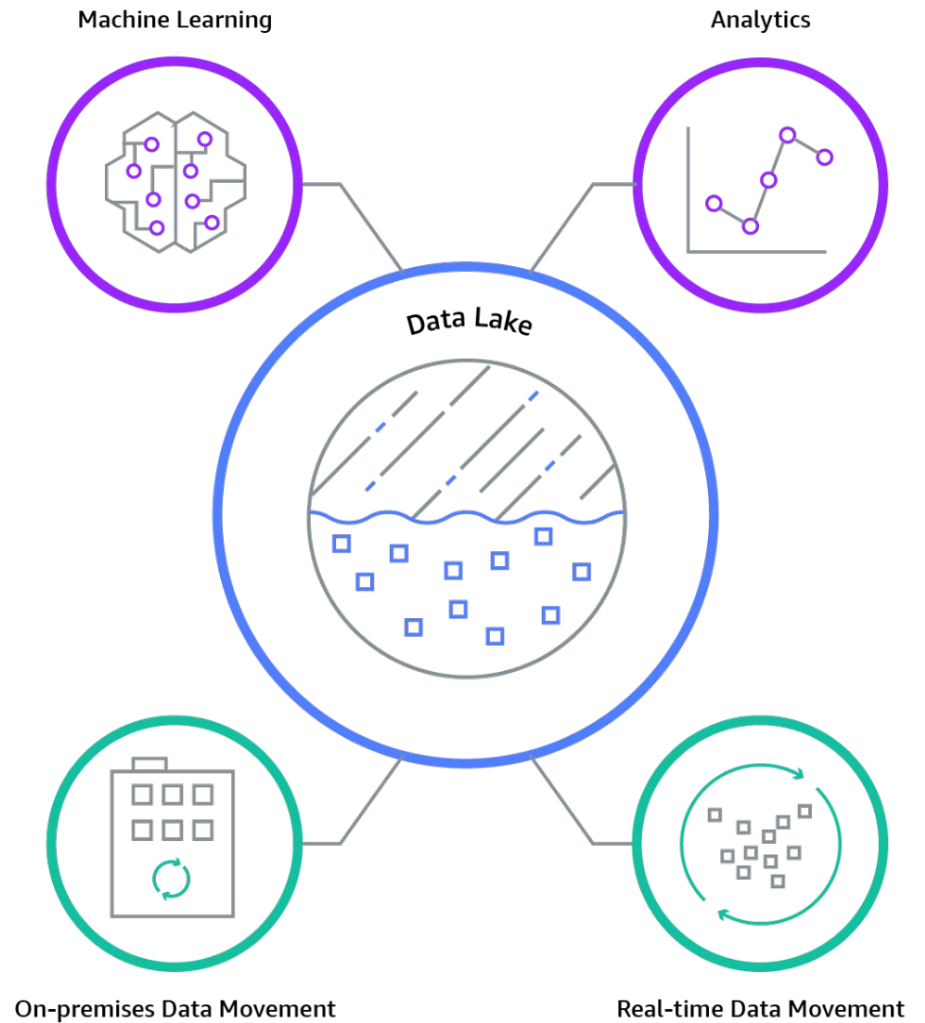


# Data Warehouse, Data Lake y Data Lakehouse: ¿cuál utilizar?





# Data Lake

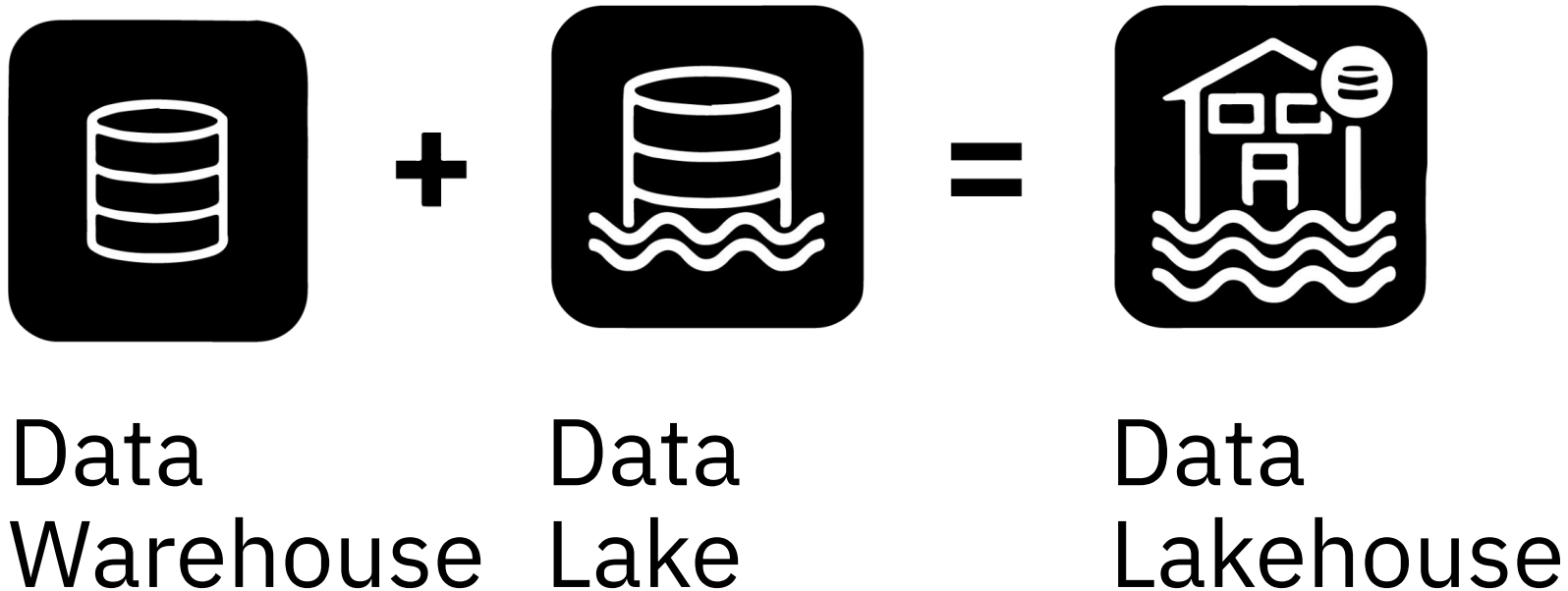


# Data Warehouse vs. Data Lake

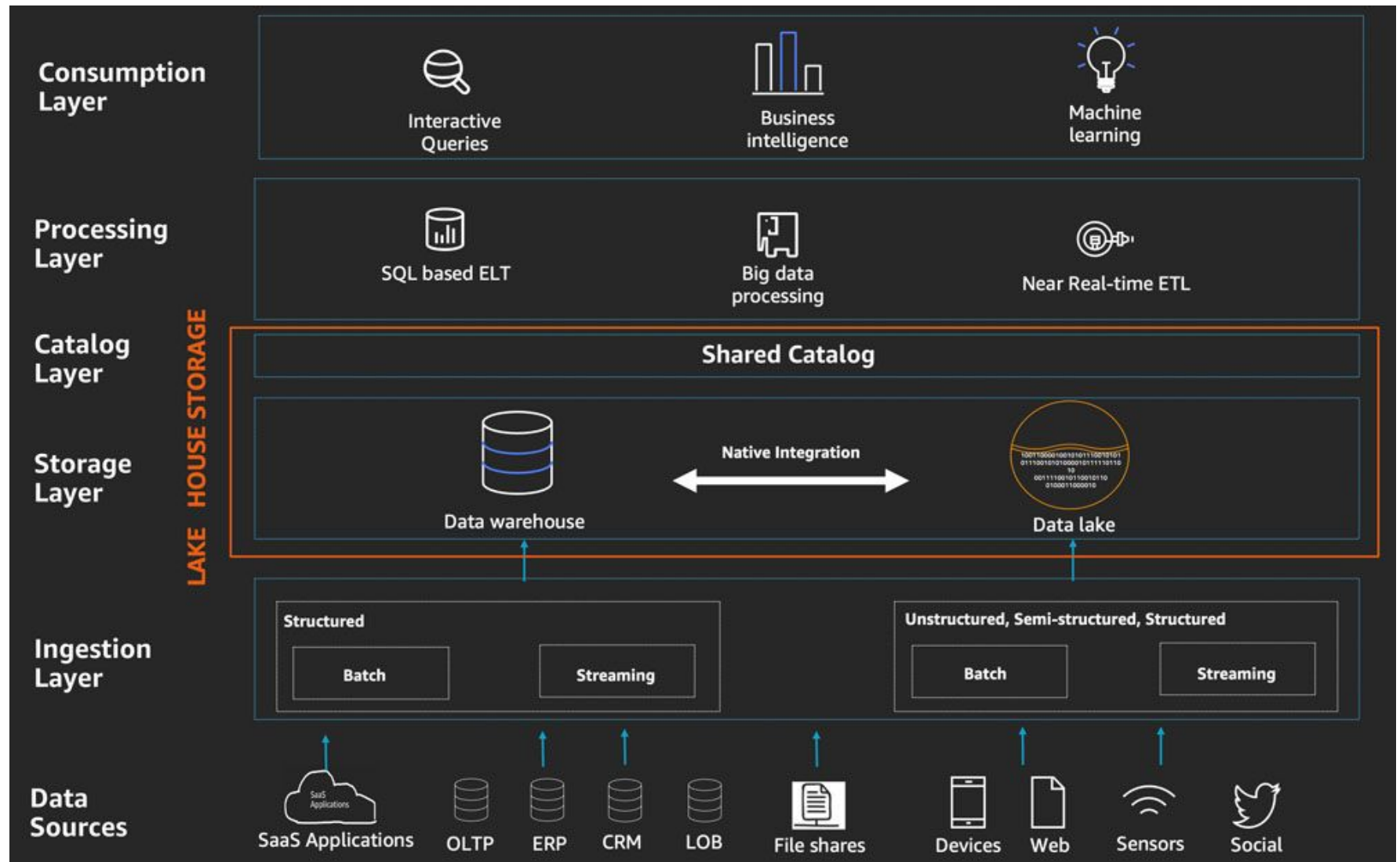
Características	Data Warehouse	Data Lake
<b>Data</b>	Optimizado para analizar datos relacionados de fuentes como BD transaccionales, operativas y aplicaciones de negocio.	Datos no relacionados de fuentes como web sites, redes sociales, dispositivos IoT, aplicaciones móviles.
<b>Schema</b>	La estructura de datos se define antes de la implementación para optimizar las consultas (schema-on-write).	Almacena información sin la definición de una estructura de datos. Permite implementar sin conocer aún las preguntas de negocio (schema-on-read).
<b>Data Quality</b>	Los datos se limpian, enriquecen y transforman para que puedan actuar como la "única fuente de verdad".	Cualquier dato que pudo o no pasar por un proceso de limpieza y transformación ( raw data).
<b>Users</b>	Analistas de negocio.	Científicos de datos, ingenieros de datos, y analistas de datos (usando información limpia).
<b>Analytics</b>	Reportes, tableros de control y BI.	Machine Learning, análisis predictivos, data discovery y profiling.

**¿Data Lakehouse?**

# Data Lakehouse



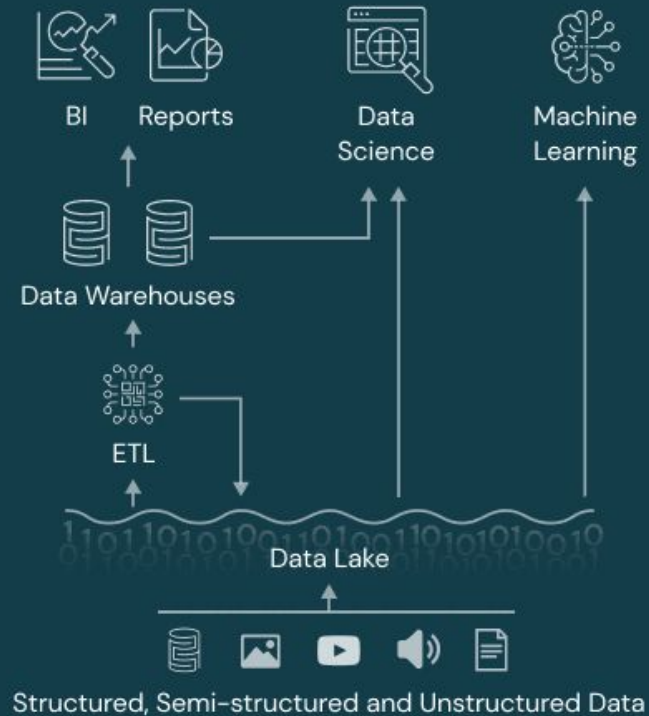
# Data Lakehouse



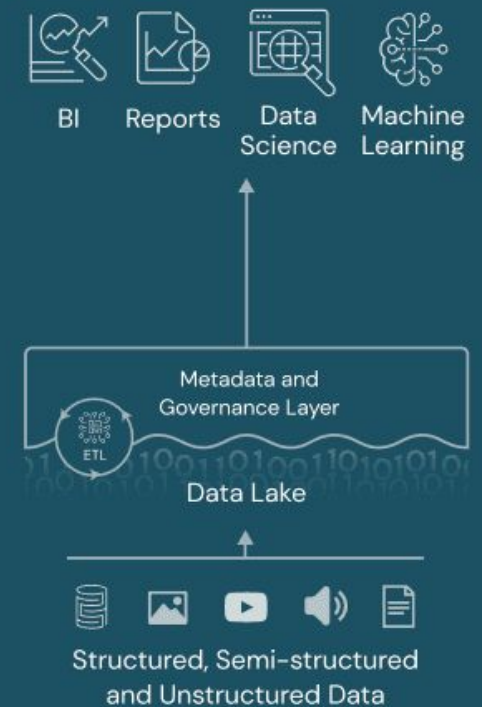
## Data Warehouse



## Data Lake



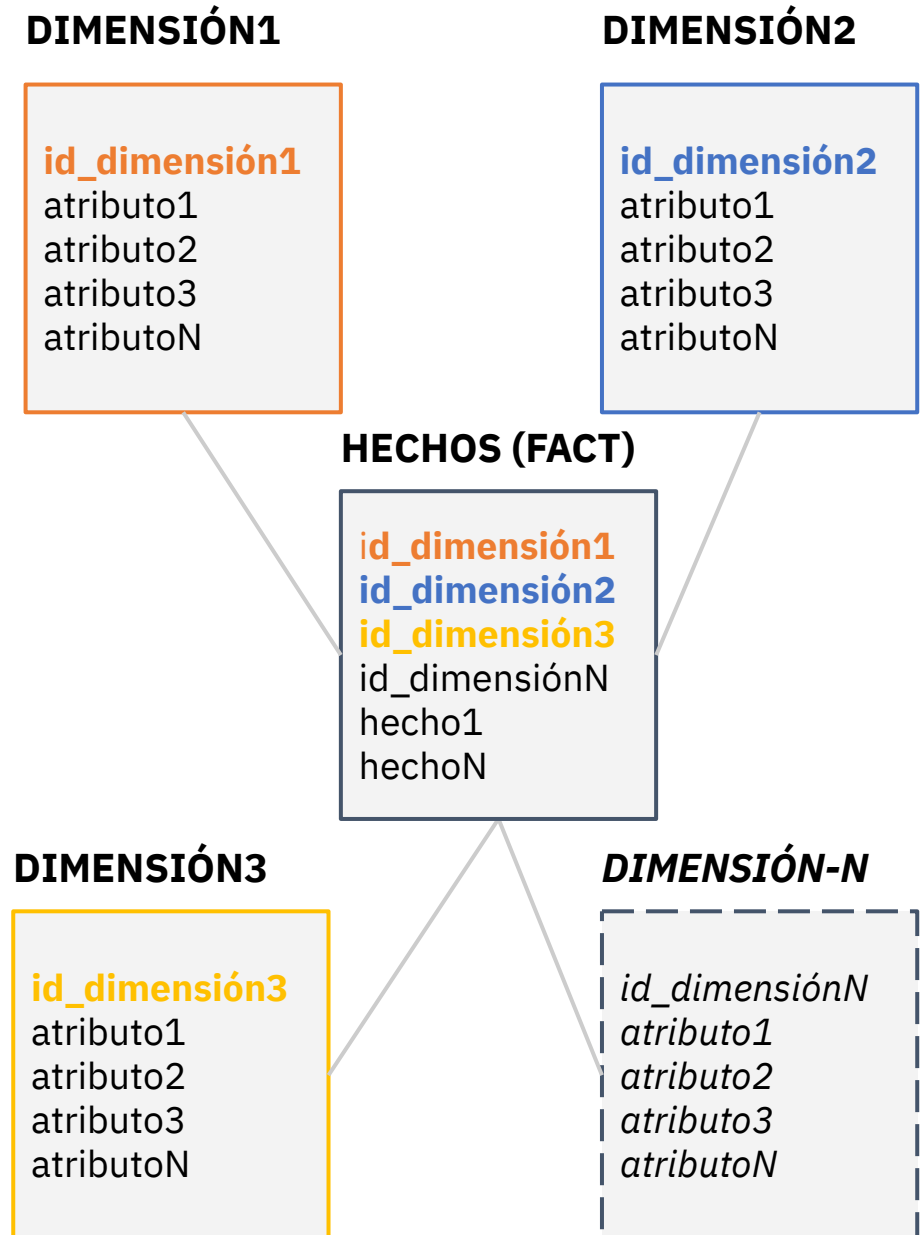
## Data Lakehouse



# Tipos de esquemas dimensionales

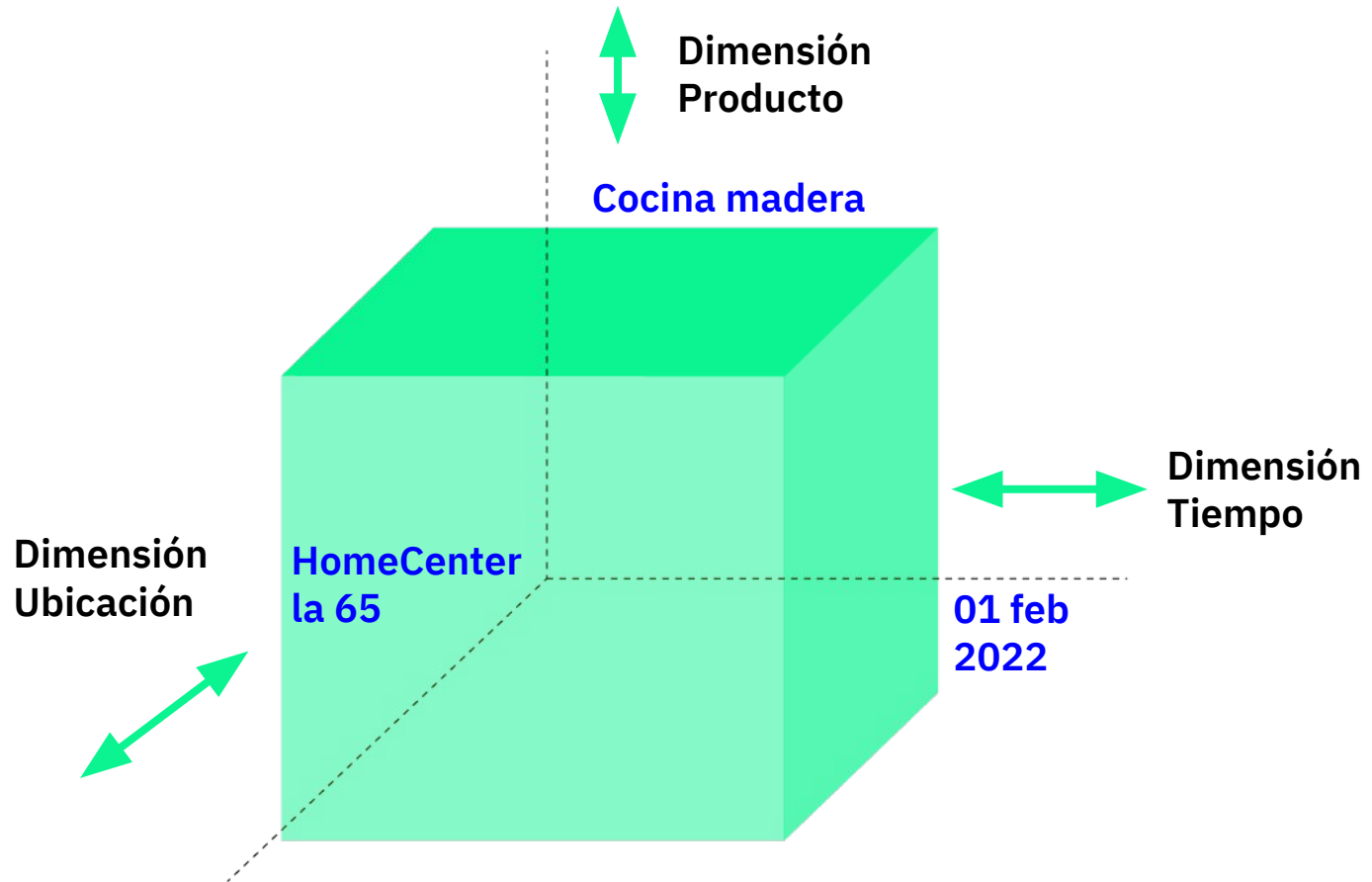
# Modelo multidimensional

Un evento objeto de análisis (**hecho**) que es evaluado bajo diferentes perspectivas o puntos de vista (**dimensiones**).



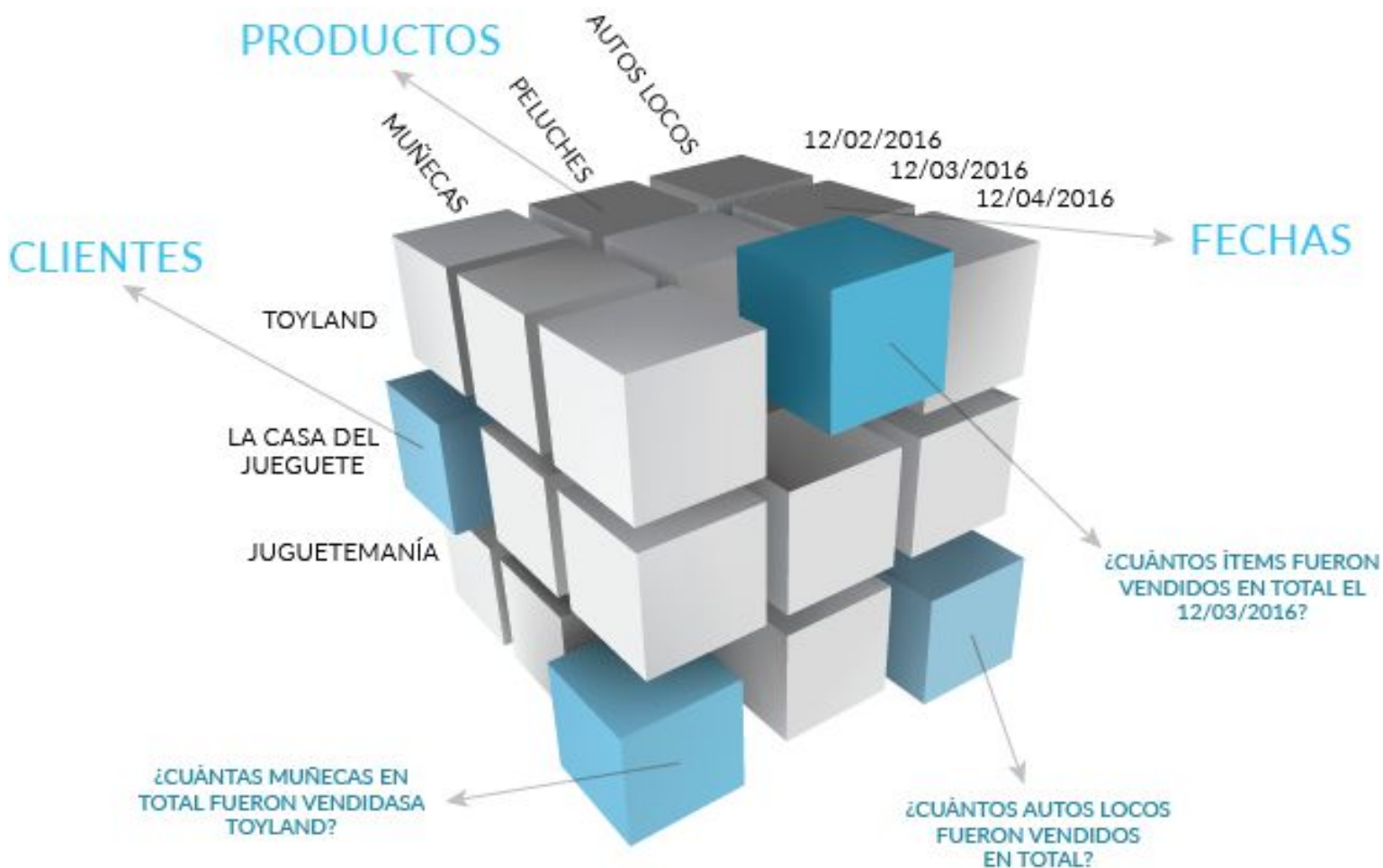


# Cubo



Conceptualmente, el modelo de DWH se puede representar mediante un cubo.

# Cubo



# Esquema estrella

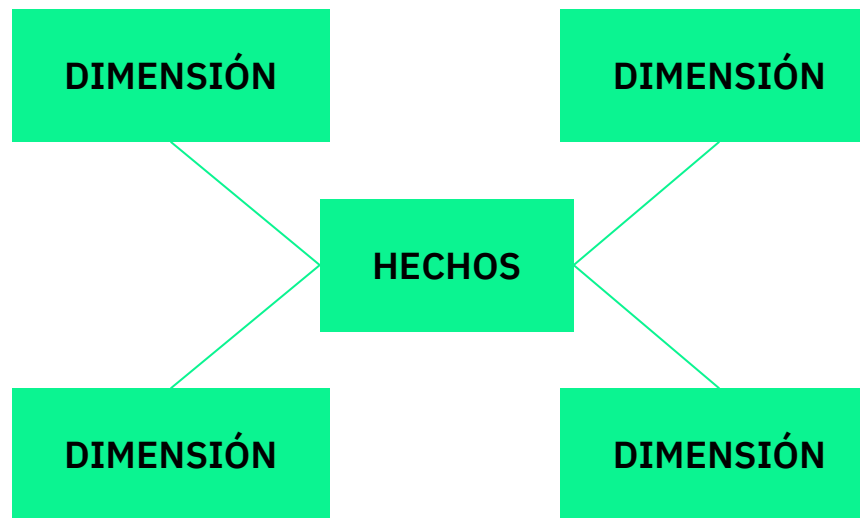
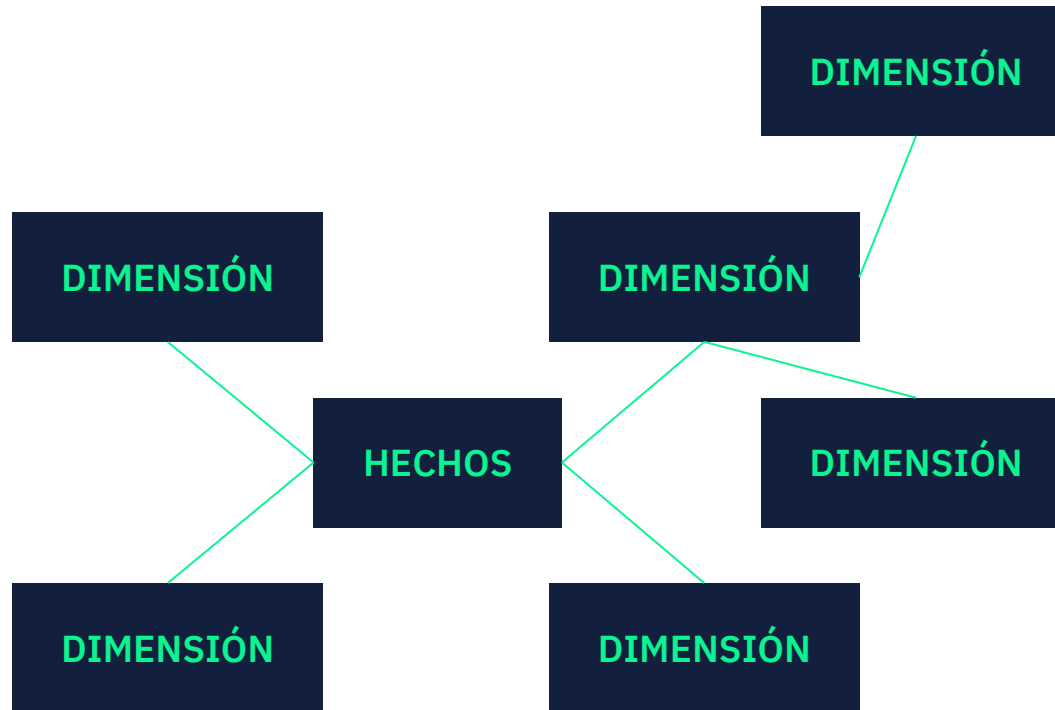


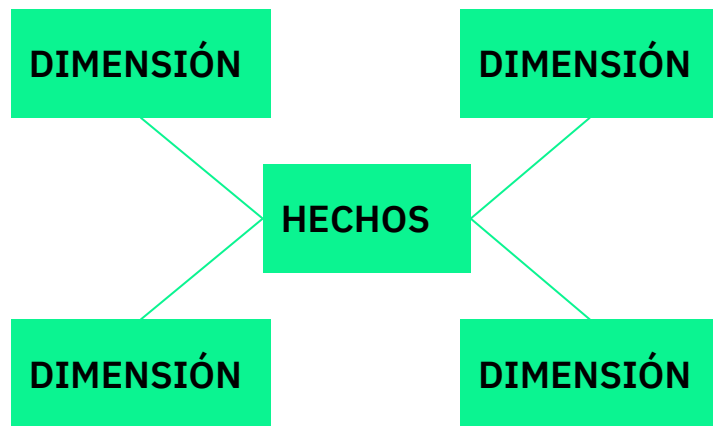
Tabla central (hechos) que se relaciona con dimensiones

# Esquema copo de nieve

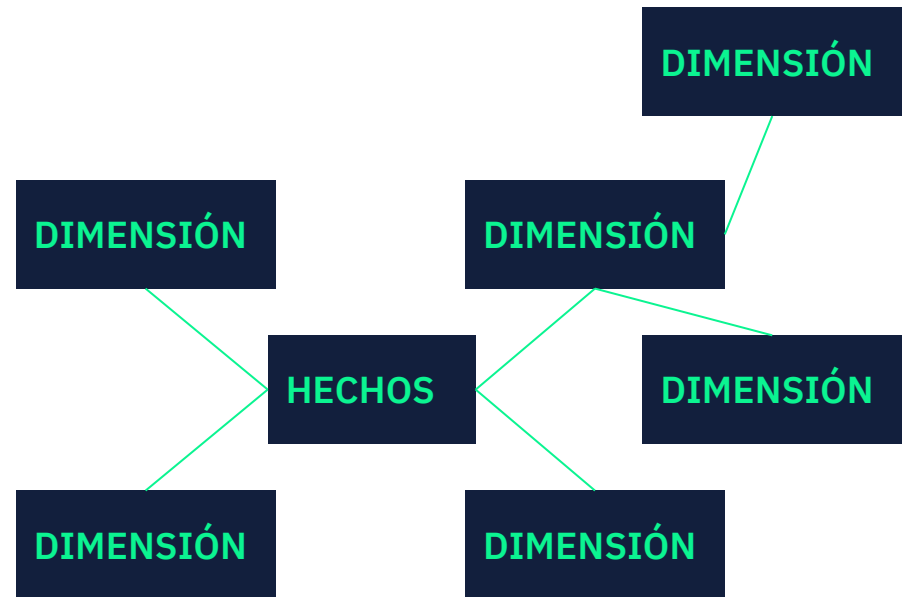


Dimensiones normalizadas que se relacionan con otras dimensiones

# Tipos de esquemas



Estrella



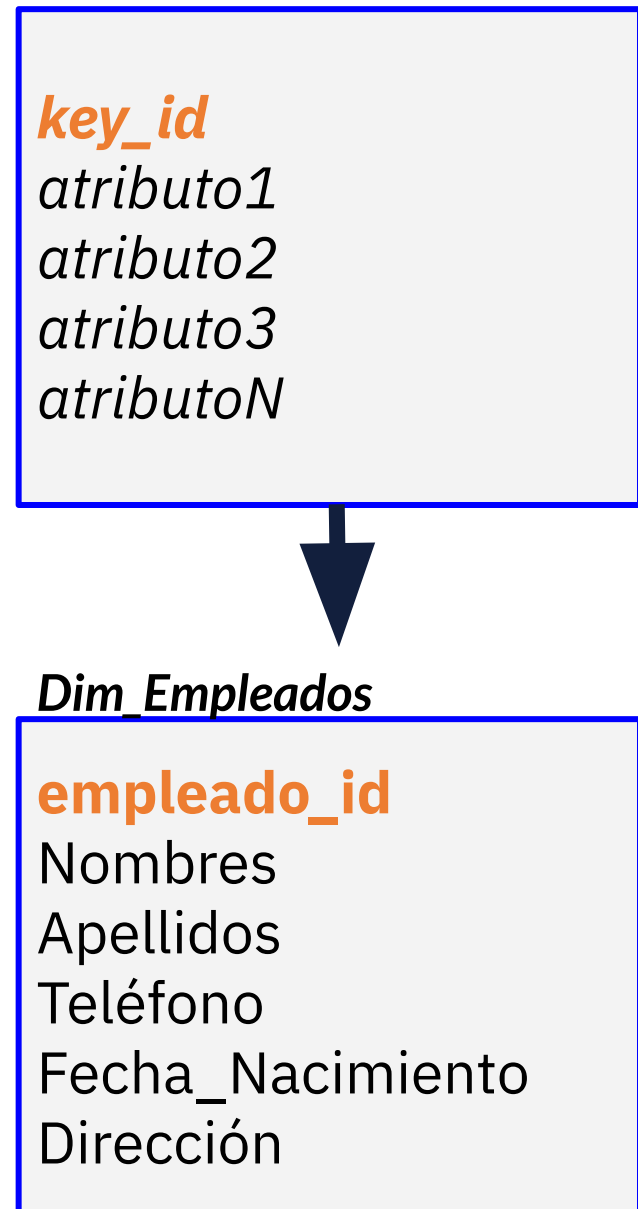
Copo de nieve

**Dimensiones  
lentamente  
cambiantes**

# Dimensión

Los diferentes actores en los procesos del negocio.

- ¿Quién?
- ¿Qué?
- ¿Cómo?
- ¿Cuándo?
- ¿Dónde?





# Atributos

- **Jerárquicos**

Permiten ir de lo general a lo particular. Consolidar y desagregar. Por ej: país.

- **Descriptivos**

Información relevante. Netamente descriptiva. Por ej: dirección, teléfono, talla, clima.



# **Atributos**

- **De control**

Datos de auditoría. No pertenecen al conocimiento del negocio. Por ej: fecha de carga.

# ↗ Tipos de SDC

- **Tipo 1**

Sobreescribir el atributo actualizado.

- **Tipo 2**

Agrega un nuevo registro con el cambio.

- **Tipo 3**

Agrega un nuevo atributo “anterior”.

**Dimensión lentamente  
cambiante tipo 1**

# SCD tipo 1 - Reemplaza

Id_Estudiante	Nombre Completo	Facultad
EST12345	Pepito Perez	Mercadeo

Cambio de facultad a **Ingeniería**

Id_Estudiante	Cod_Estudiante	Nombre Completo	Facultad
1	EST12345	Pepito Perez	Ingeniería

**Dimensión lentamente  
cambiante tipo 2**

# SCD tipo 2 - Agrega fila

Id_Estudiante	Nombre Completo	Facultad
EST12345	Pepito Perez	Mercadeo

Cambio de facultad a **Ingeniería**

Id_Estudiante	Cod_Estudiante	Nombre Completo	Facultad	Start_date	End_date
1	EST12345	Pepito Perez	Mercadeo	01/01/2020	01/01/2023
2	EST12345	Pepito Perez	Ingeniería	02/01/2023	31/12/9999

**Dimensión lentamente  
cambiante tipo 3**

# SCD tipo 3 - Agrega atributo

Id_Estudiante	Nombre Completo	Facultad
EST12345	Pepito Perez	Mercadeo

Cambio de facultad a **Ingeniería**

Id_Estudiante	Cod_Estudiante	Nombre Completo	Facultad_old	Facultad_new
1	EST12345	Pepito Perez	Mercadeo	Ingeniería

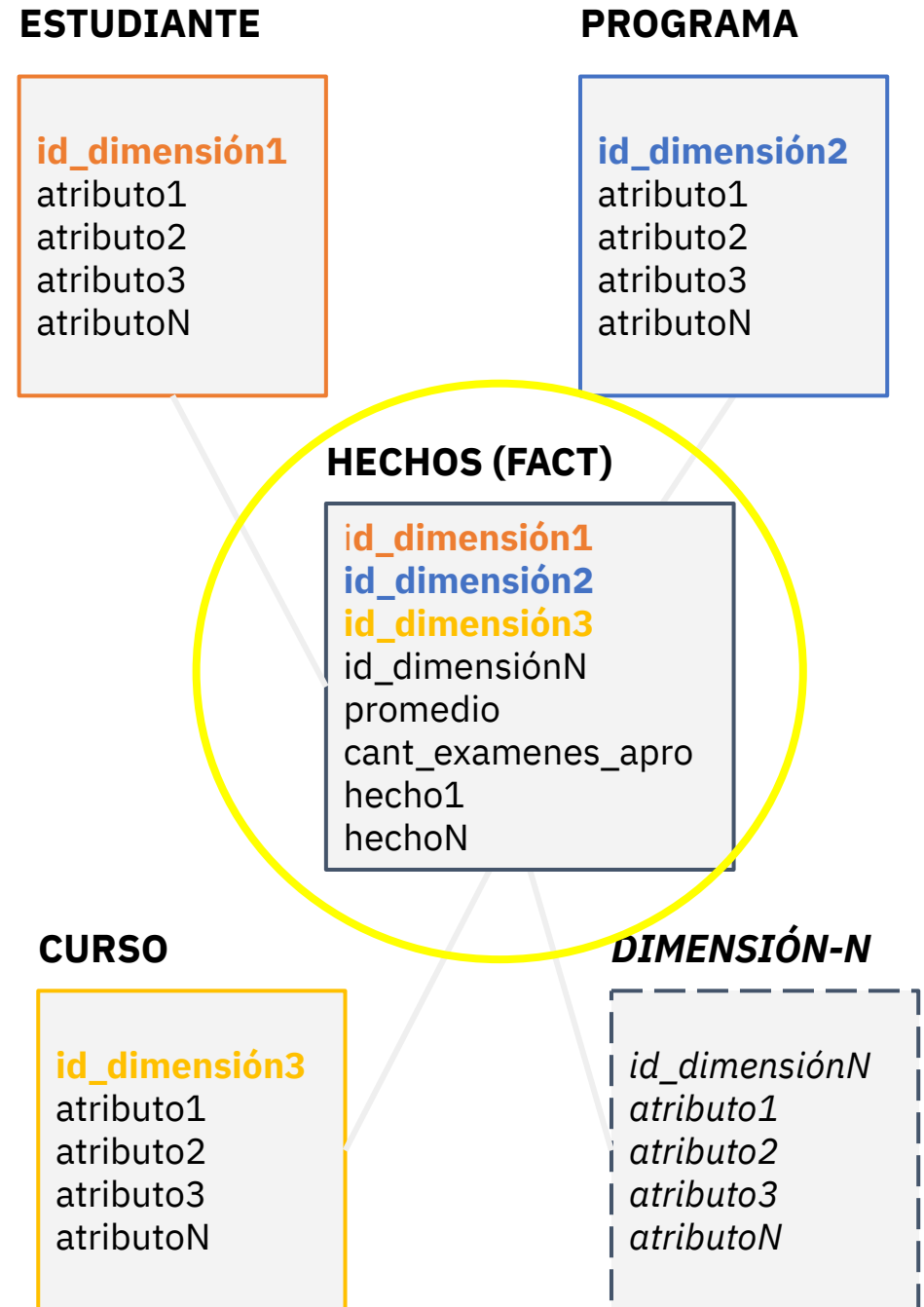


**Solución reto anterior**

# Tabla de hechos

# Hechos (Fact)

- Contienen información cuantitativa de un proceso de negocio.
- **Medidas - Métricas.**
- Contiene claves foráneas (llaves de las dimensiones).



# Configuración del Setup

- Postgres:
- Redshift:
- Pentaho:
- BD AdventureWorks:

[https://wiki.postgresql.org/wiki/Sample\\_Databases](https://wiki.postgresql.org/wiki/Sample_Databases)

# Identificación de dimensiones y métricas

# Preguntas del negocio

“

**El trabajo consiste en hacer preguntas, todas las que se puedan, y hacer frente a la falta de respuestas precisas con una cierta humildad.**

”

*Arthur Miller*



# Preguntas de Negocio

- **Unidades vendidas** de cada **producto** por **cliente** en un **tiempo** determinado.
- **Cantidad de contrataciones** por **área** en un **país** específico.

# Ejemplo

## Organización:

Cadena de supermercados.

## Actividad objeto de análisis:

Ventas de productos.

## Información registrada sobre una venta:

Del **producto** “crema dental” se han vendido en el **almacén** “Almacén nro.1” el **día** 2/2/2030, 5 **unidades** por un **valor** de \$20.

# Ejercicio

- ¿**Cuánto** ha sido en cantidades y valores, los descuentos y las ventas netas (venta-descuento), para cada mes y día?
- ¿**Cuánto** ha crecido o disminuido las ventas netas al corte del mes de Marzo del año 2013 para cada vendedor?
- ¿**Cuál** es el producto más vendido a corte del día? Por categoría.
- ¿**Quién** es el cliente que más unidades ha comprado en el último año?

- ¿**Cuánto** ha sido en cantidades y valores, los descuentos y las ventas netas (venta-descuento), para cada mes y día?
- ¿**Cuánto** ha crecido o disminuido las ventas netas al corte del mes de Marzo del año 2013 para cada vendedor?
- ¿**Cuál** es el producto más vendido a corte del día? Por categoría.
- ¿**Quién** es el cliente que más unidades ha comprado en el último año?

**Tu reto**

- ¿**Dónde** geográficamente se ubican las mayores ventas netas del mes de Diciembre del año 2014?
- ¿**Quién** es el vendedor que más ventas realizó por semestre, teniendo en cuenta el nivel del cargo en el momento de la venta?

- ¿**Dónde** geográficamente, se ubican las mayores **ventas netas** del **mes** de Diciembre del **año** 2014?
- ¿**Quién** es el **vendedor** que más **ventas** realizó por **semestre**, teniendo en cuenta el **nivel del cargo** en el momento de la venta?



# Diseño de modelo

# Reglas de negocio

- Crear un campo con el nombre completo del cliente.
- El campo de observación del producto es demasiado largo. Recortar a los 100 primeros caracteres.
- Si un vendedor tiene personas a cargo, marcarlo como beneficiario del bono.

# Ejercicio

# Documento de mapeo

ETL

**Identificar  
transformaciones**

v, h

1, 0

varón,  
hembra



V, H

dd/mm/yyyy

yyyy-mm-dd

mm-yyyy



dd/mm/yyyy





**Código de producto = 57B1050**

**código  
país**

**zona de  
inventario**

**código de  
producto**



**Tu reto**

# Creación del modelo físico

ETL

**Tu reto**

# Extracción

ETL - Dimensiones

# Transformación

ETL - Dimensiones

# Carga

ETL - Dimensiones

**Tu reto**

# Extracción

ETL - Hechos



# Transformación

ETL - Hechos

# Carga

ETL - Hechos

# Orquestar ETL

# Ejercicio para el estudiante

# Reflexiones y cierre

- Repaso
- Dificultades
- Siguietes pasos: Curso de BI con Power BI
- Enviar el proyecto
- Autoevaluarse en los comentarios
- Evaluar este curso
- Redes
- Tchau