



Escola Politécnica de Pernambuco

Especialização em Ciência de Dados e Analytics

Introdução à Ciência de Dados

Aula 3

Prof. Dr. Alexandre Maciel
alexandre.maciel@upe.br

QUAIS DADOS ANALISAR?



PRODUÇÃO POR INPUT

- Teclado.
- Mouse.
- *Touch Screen.*
- Scanners.
- Código de Barras.
- RFID.
- Câmeras.
- Filmadoras.



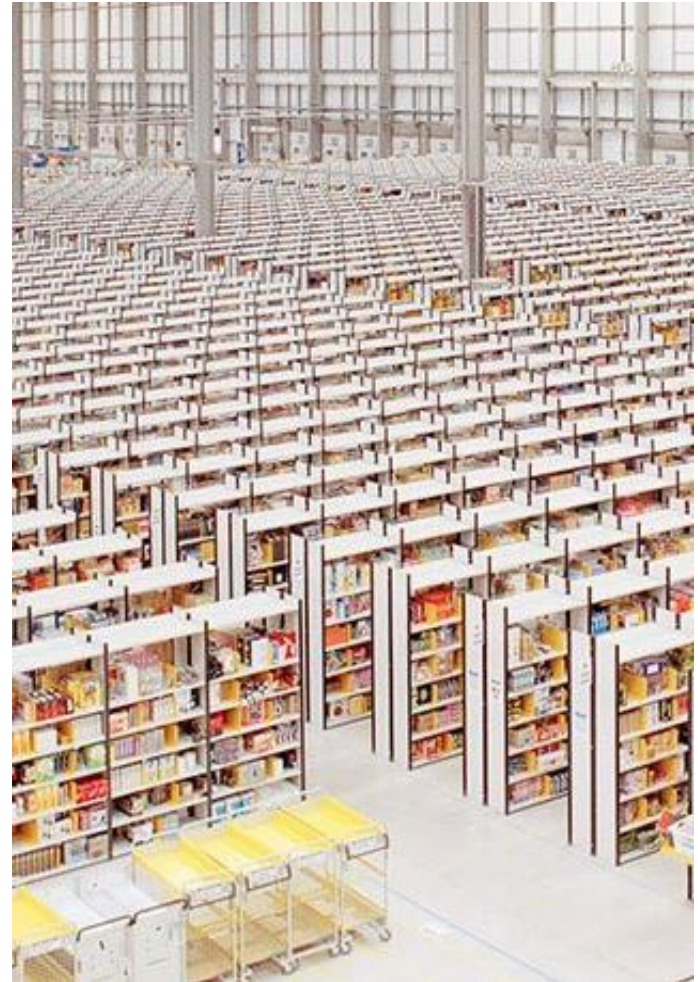
DADOS POR PROCESSAMENTO

- **Análise ou Execução de Procedimentos.**
- **Criação de Modelos:**
 - Estatísticos;
 - Machine Learning.
- **Criação de Data Warehouses.**
 - Data Marts.



DADOS POR TRANSFORMAÇÃO

- **Data Warehouse**
- **Objetivo:**
 - Geração de informação e conhecimento.
- **OLAP versus OLTP.**
- **Data Marts:**
 - Unidades temáticas.
- **Base para os BI's.**



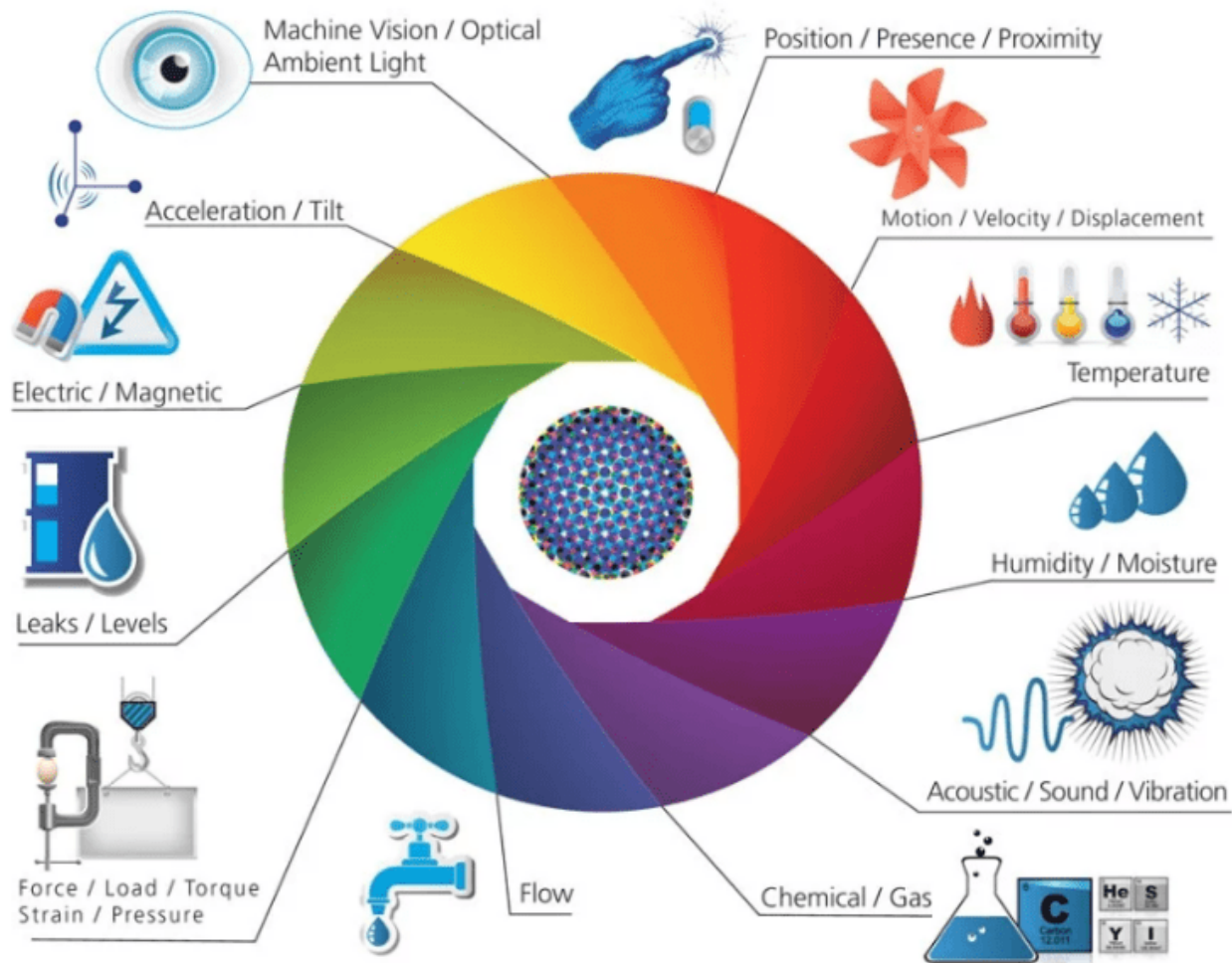
DEFINIÇÕES

Business Intelligence

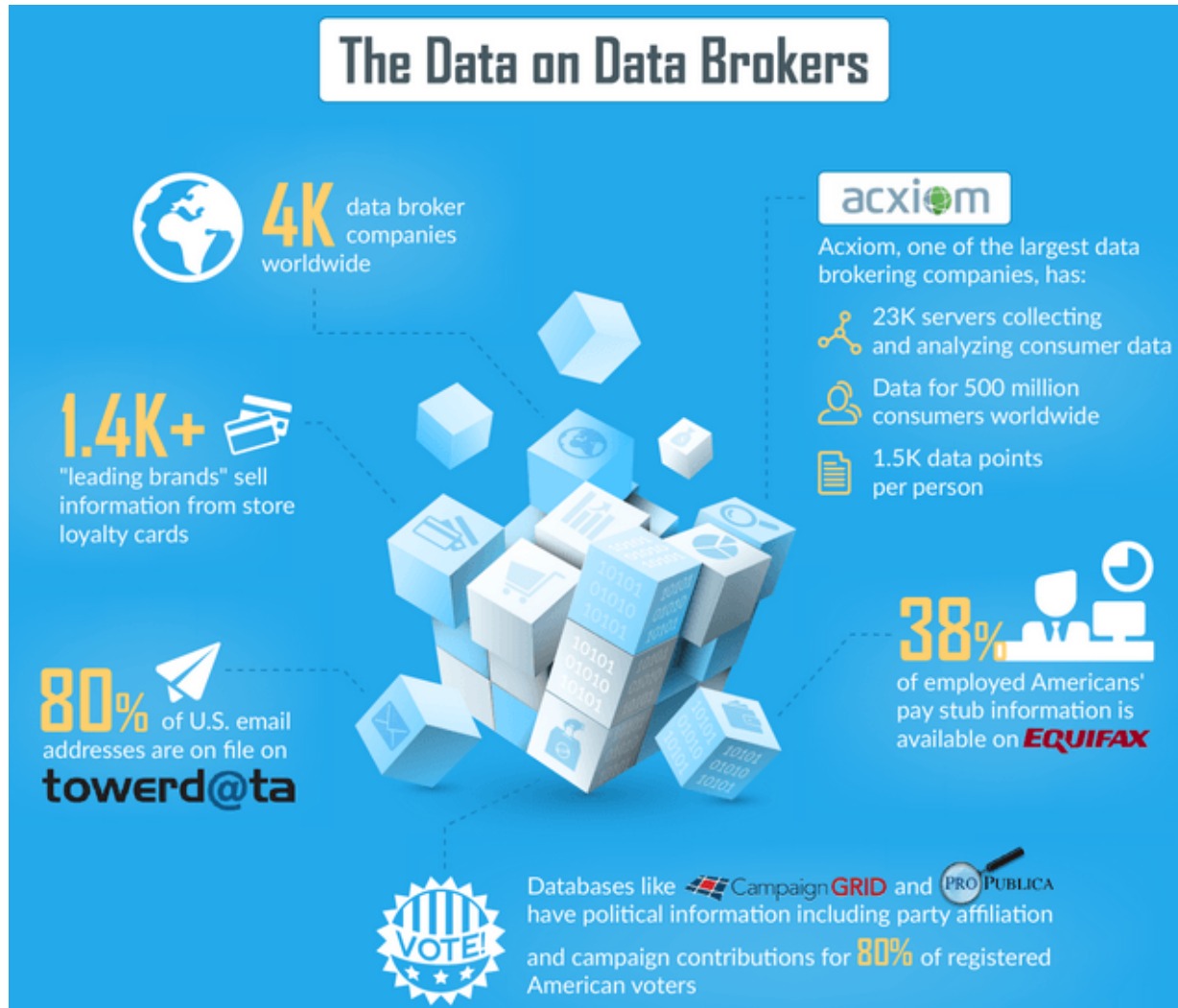
se refere à **metodologia**,
ferramentas, **técnicas** de
produzir dados para apoio
à **decisões**.
(Amaral).



PRODUÇÃO POR SENSORES

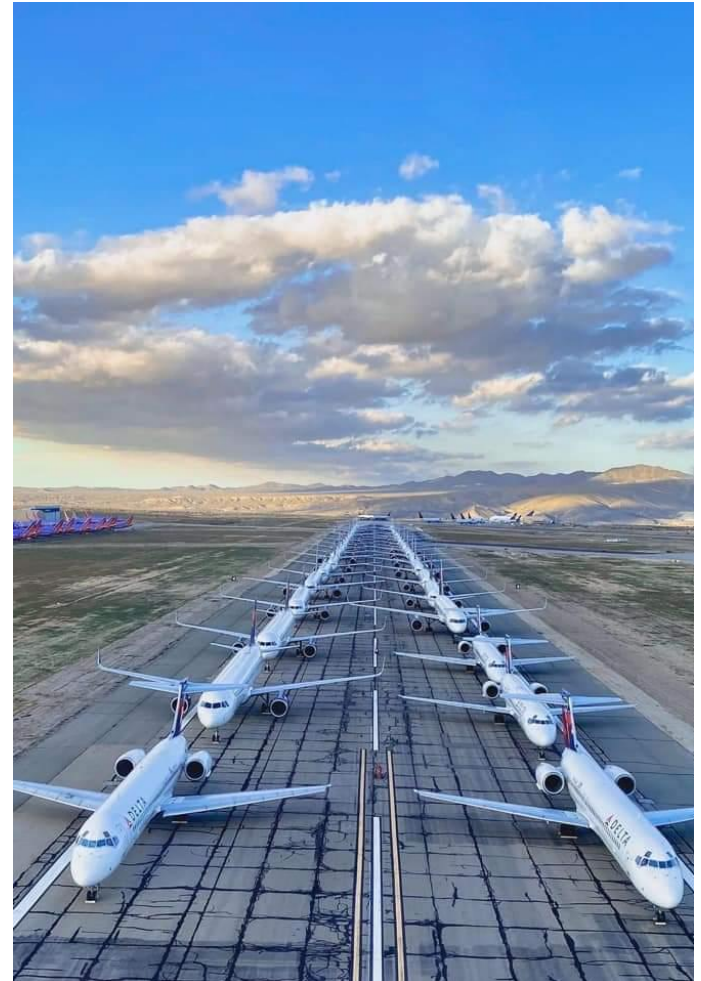


DADOS COMPRADOS

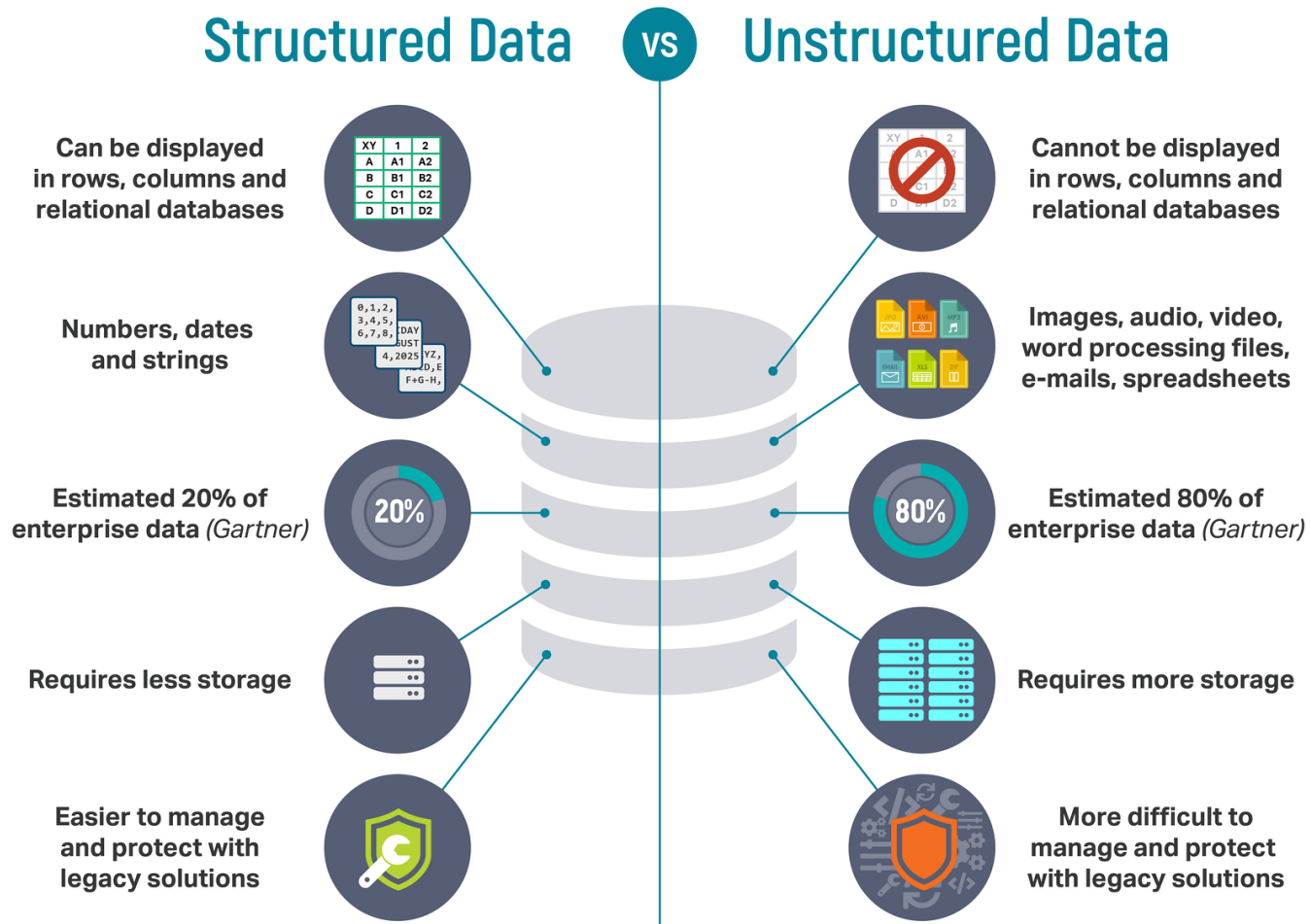


ARMAZENAMENTO

- Dado produzido, armazenado, recuperado, utilizado.
- Segurança, integridade, minimização de redundância, concorrência, otimização de espaço...



ARMAZENAMENTO



SEMI-ESTRUTURADO

XML

Choose File sample.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <article xmlns="http://docbook.org/ns/docbook"
3   xmlns:xlink="http://www.w3.org/1999
4   <info>
5     <title>Welcome to DocBook Suppo
6   </info>
7   <sect1>
8     <title>Inline Markup and Images
9     <para>This sample shows that &l
10       with the
11       dockbookx.dtd<?oxy_delete a
12     <para>The following <code>Docbo
13     transformation scenario. Fo
```

Download Copy Clear

→

JSON

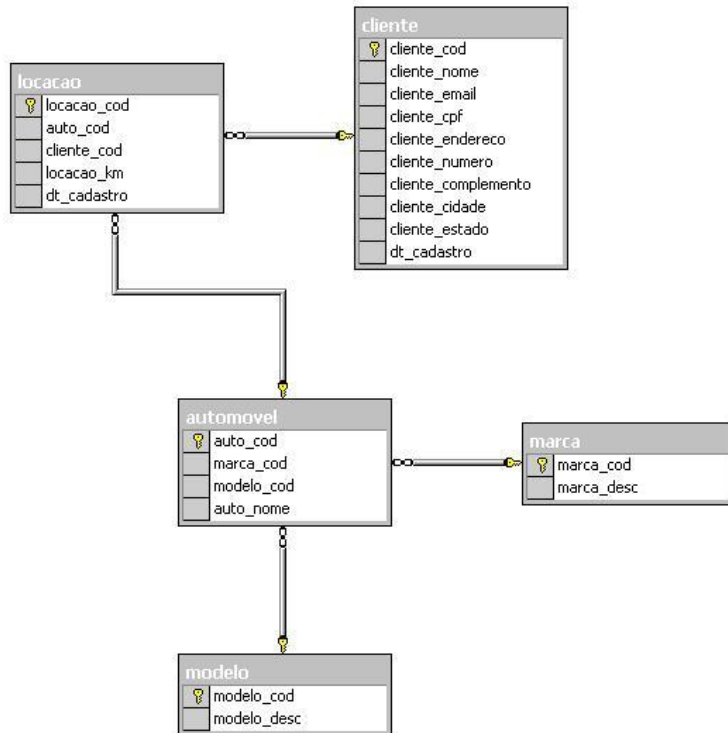
Choose File No file chosen

```
1 {
2   "article": {
3     "xmlns": "http://docbook.org/ns/docb
4     "version": "5.0",
5     "xmlns:xlink": "http://www.w3.org/19
6     "info": {"title": "Welcome to DocBoo
7     "sect1": [
8       {
9         "title": "Inline Markup and Imag
10        "para": [
11          [
12            "This sample shows that <oXy
13          {
```

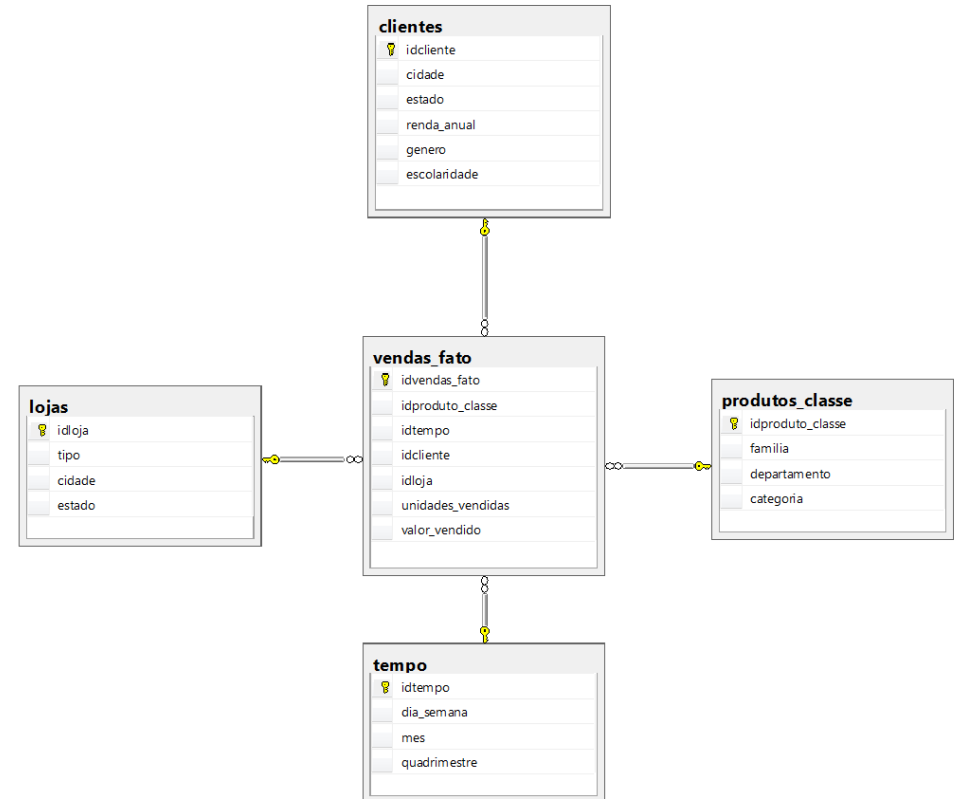
Download Copy Clear

←

MODELO RELACIONAL X DIMENSIONAL

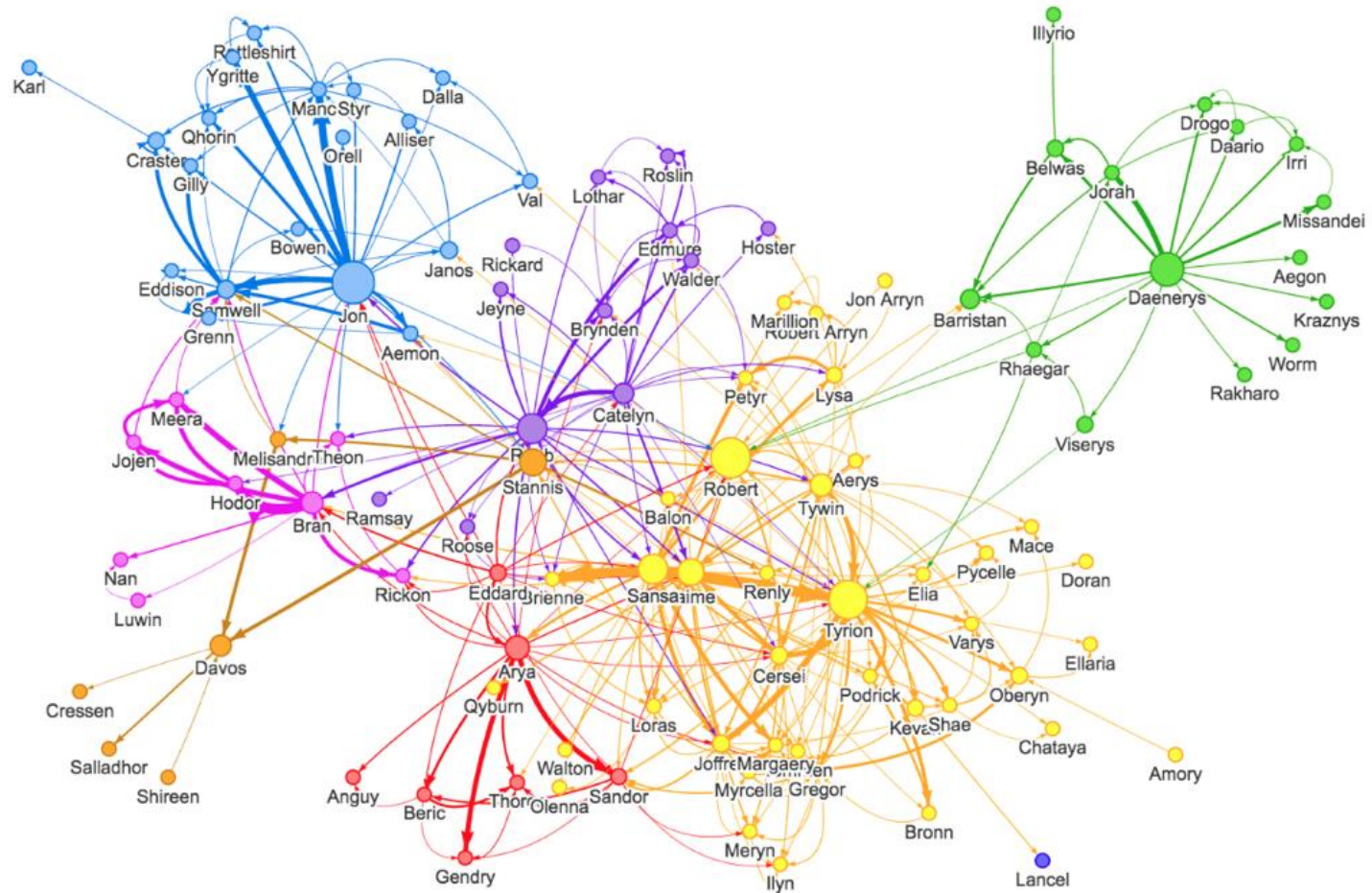


RELACIONAL



DIMENSIONAL

DADOS EM GRAFOS



DADOS ORDENADOS

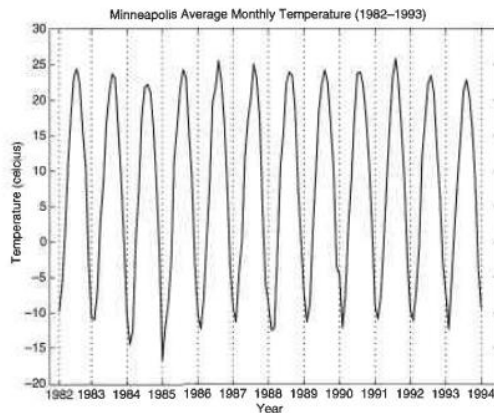
Time	Customer	Items Purchased
t1	C1	A, B
t2	C3	A, C
t2	C1	C, D
t3	C2	A, D
t4	C2	E
t5	C1	A, E

Customer	Time and Items Purchased
C1	(t1: A,B) (t2:C,D) (t5:A,E)
C2	(t3: A, D) (t4: E)
C3	(t2: A, C)

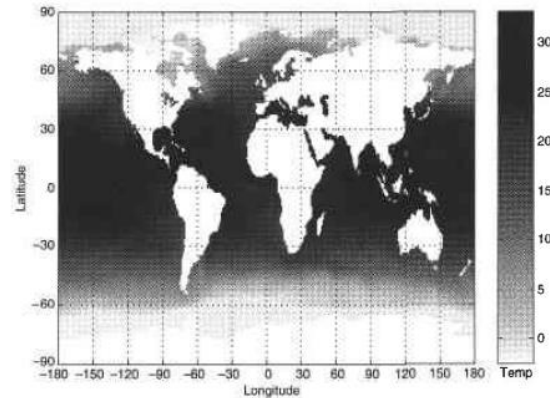
(a) Sequential transaction data.

```
GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCGG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG
```

(b) Genomic sequence data.



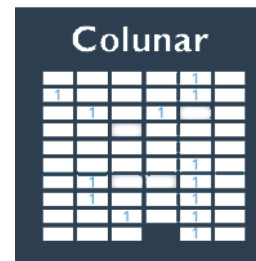
(c) Temperature time series.



(d) Spatial temperature data.

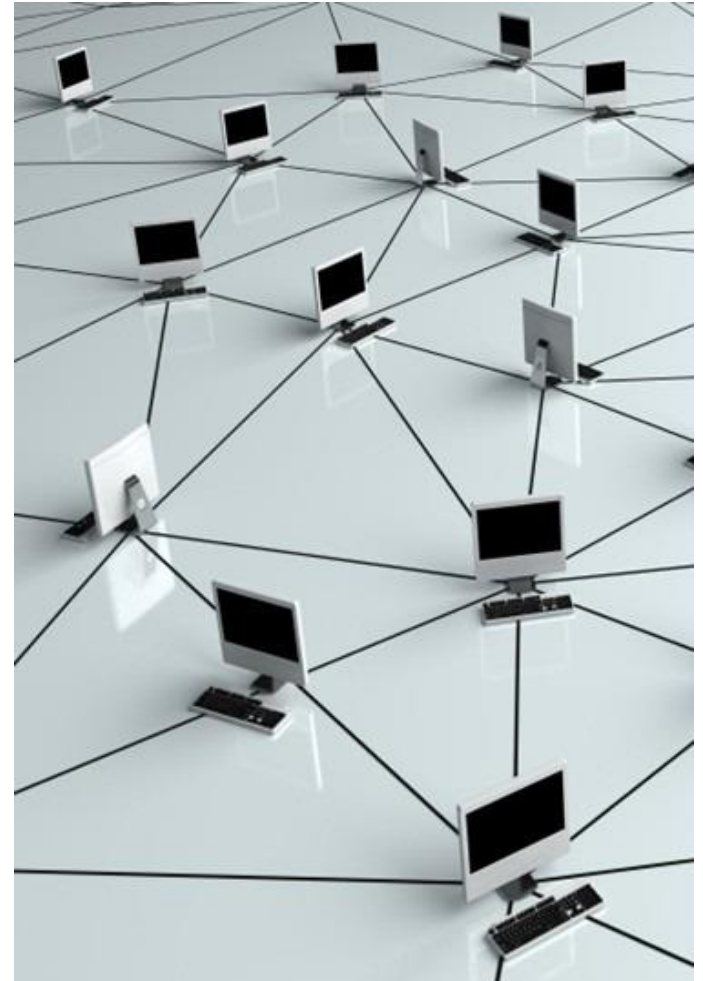
NoSQL

- Not Only SQL: Bancos de dados não relacionais.
- Agrega grandes volumes de dados.
- Famílias:



ARMAZANAMENTO DISTRIBUÍDO

- **HDFS:**
 - *Hadoop Distributed File System.*
- **Armazena grande quantidade de dados.**
- **Acesso rápido e fácil:**
 - MapReduce.
- **Tolerância a falhas.**



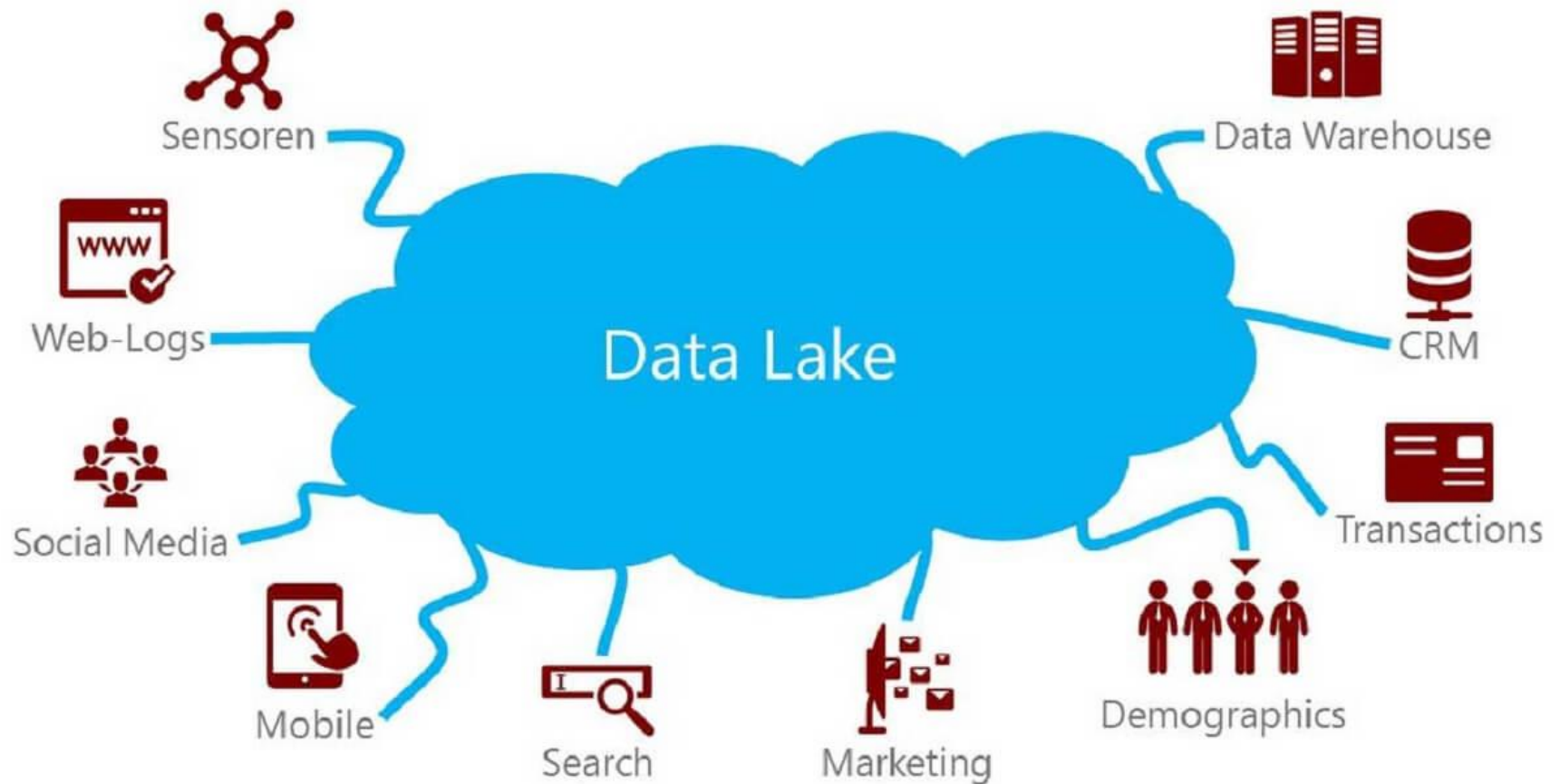
DEFINIÇÕES

Data Lake

repositório centralizado que permite armazenar todos os seus dados **estruturados** e **não estruturados** em **qualquer escala**.
(AWS).



DATA LAKE

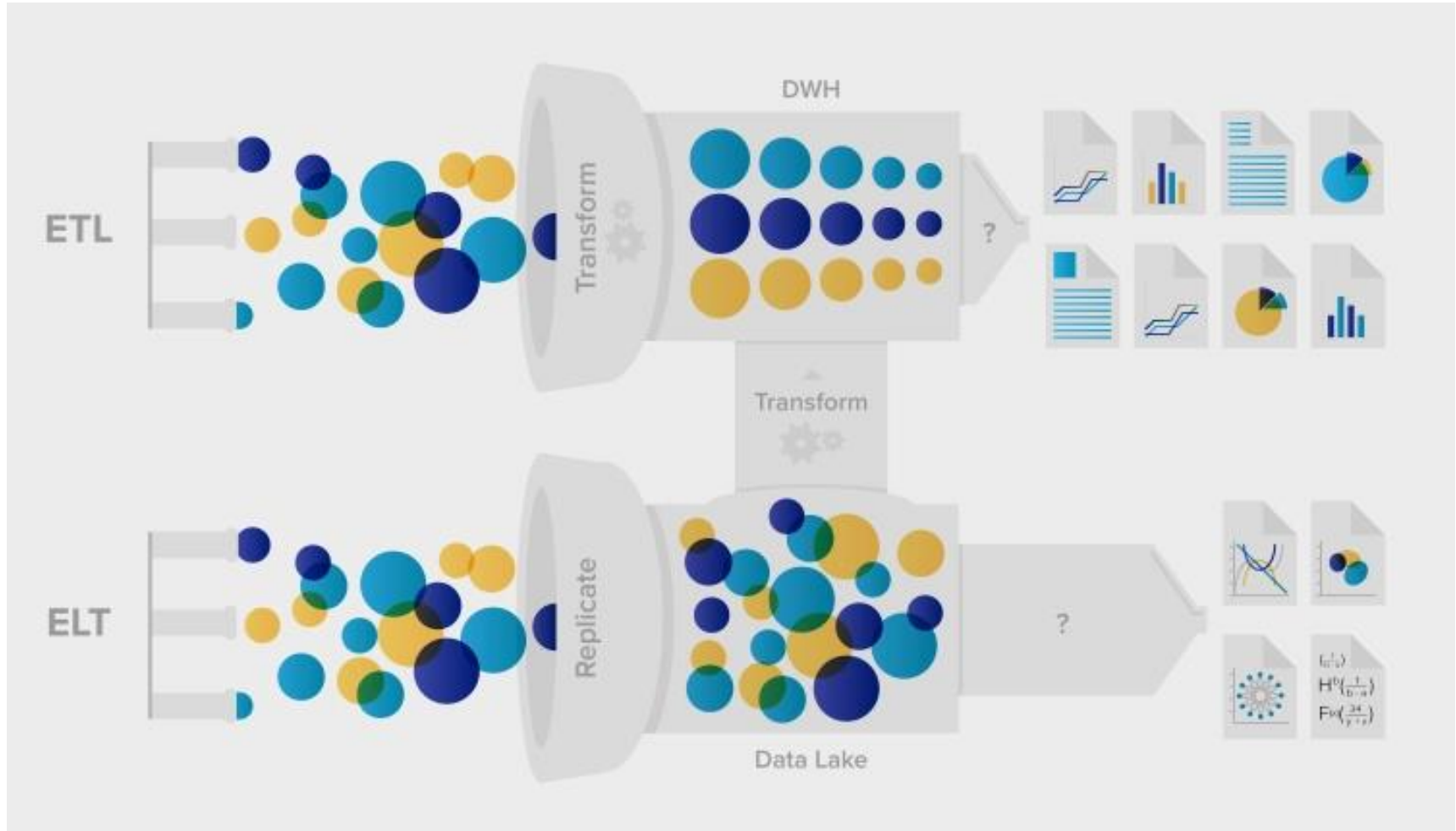


Key Differences Between the Data Lake and Data Warehouse

DATA WAREHOUSE	vs.	DATA LAKE
Structured, processed	DATA	Structured/semi-structured/unstructured/raw
Schema-on-write	PROCESSING	Schema-on-read
Expensive for large data volumes	STORAGE	Designed for low-cost storage
Less agile, fixed configuration	AGILITY	Highly agile, configure and reconfigure as needed
Mature	SECURITY	Maturing
Business pros	USERS	Data scientists et al.

Analysis Source: "A Big Data Cheat Sheet: What Marketers Want to Know" by Tamara Dull

TRANSFORMAÇÃO DOS DADOS



ATIVIDADE 3

1. **Apresente o dicionário de dados do seu projeto.**
2. **Descreva em mais detalhes os principais atributos que serão utilizados no seu projeto.**
3. **Entregue os slides (+2) de sua apresentação.**
4. **Faça o upload no Google Classroom.**