

# Business { [Intelligence]

Prof. Juscelino Fernandes da Costa Junior

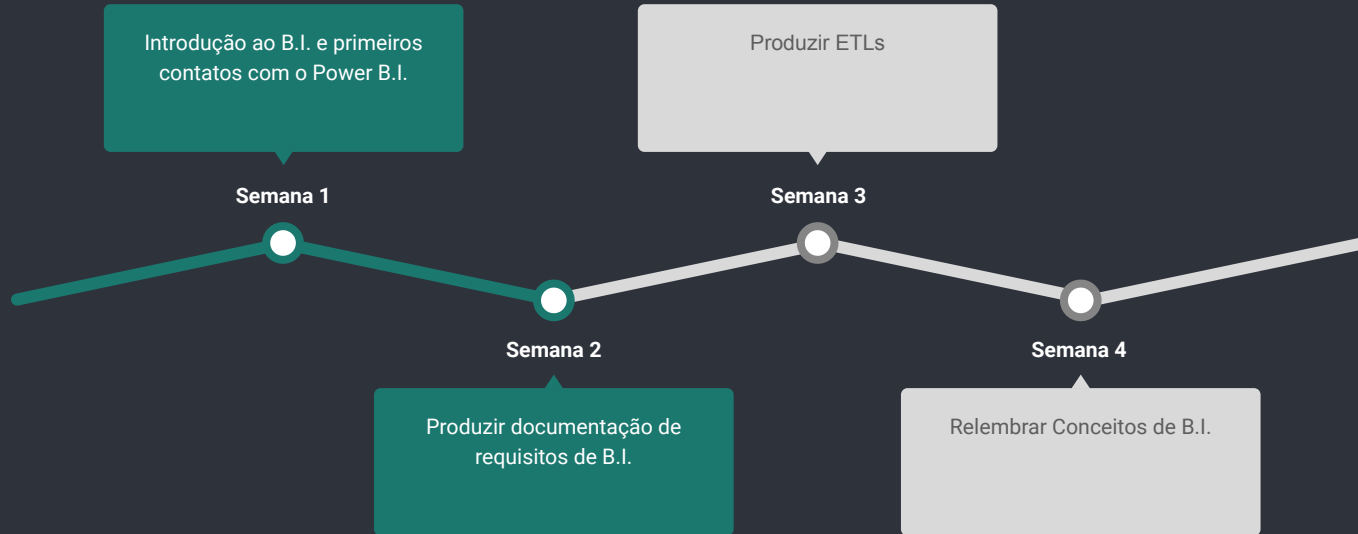
}

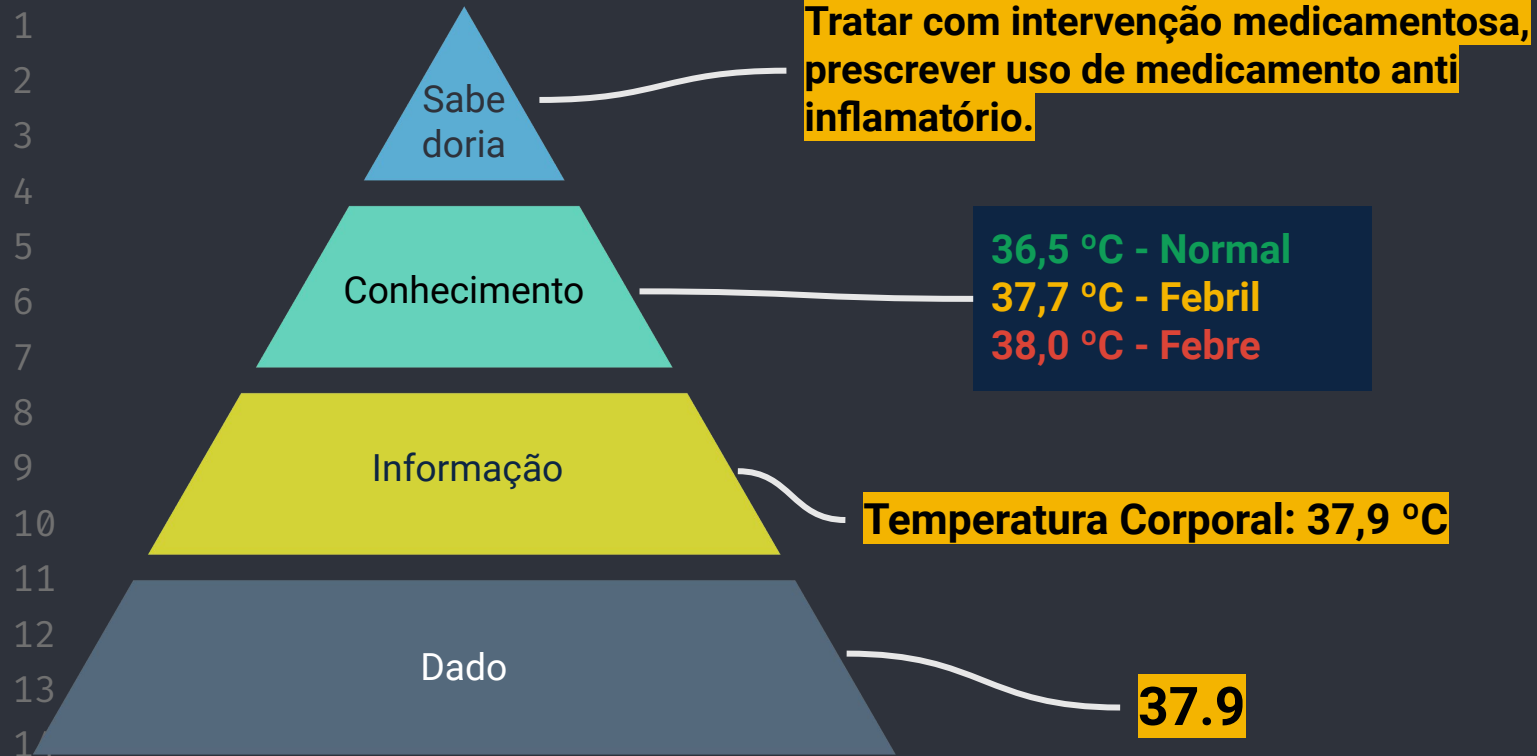
# Business Intelligence;



< B.I. é o conjunto de processos e ferramentas que proporcionam o processamento e transformação de dados brutos em informações valiosas para a tomada de decisão e formação de conhecimento >

# Business Intelligence;





# Papéis na Área de Dados {



## Engenheiro de Dados

< Computação em Nuvem,  
Big Data, Spark,  
Airflow, SQL >



## Report Designer

< Figma, Tableau,  
Superset >



## Analista de Dados

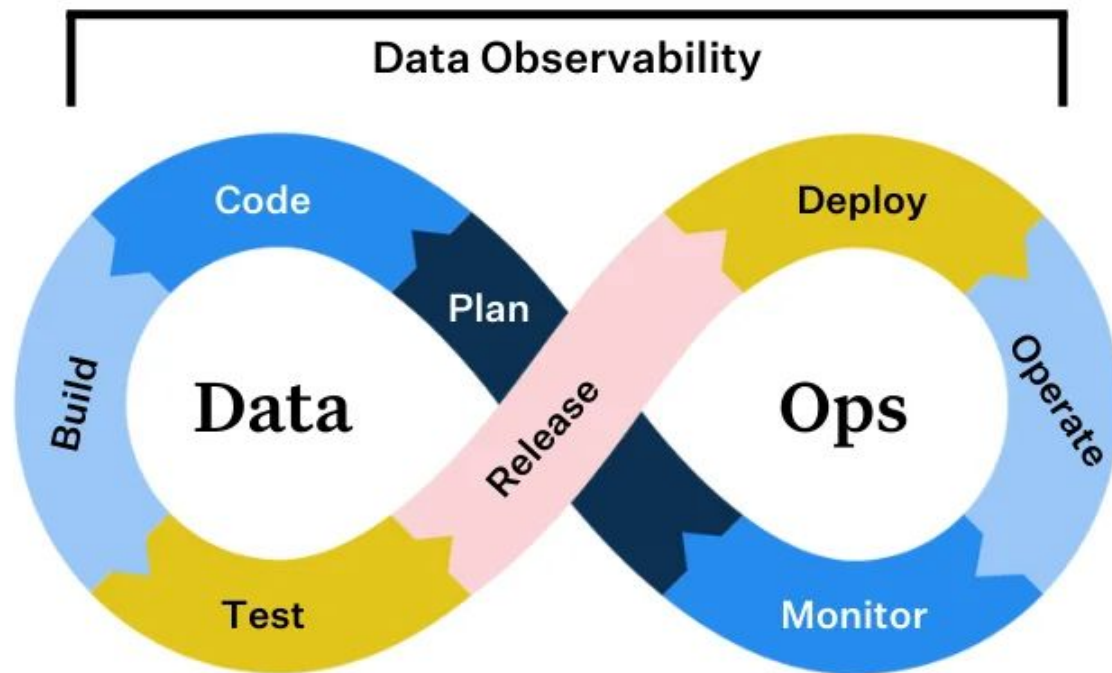
< Python, SQL, PowerBi  
Tableau, Qlickview, Excel>



## Cientista de Dados

< Python, Estatística,  
Machine Learning >

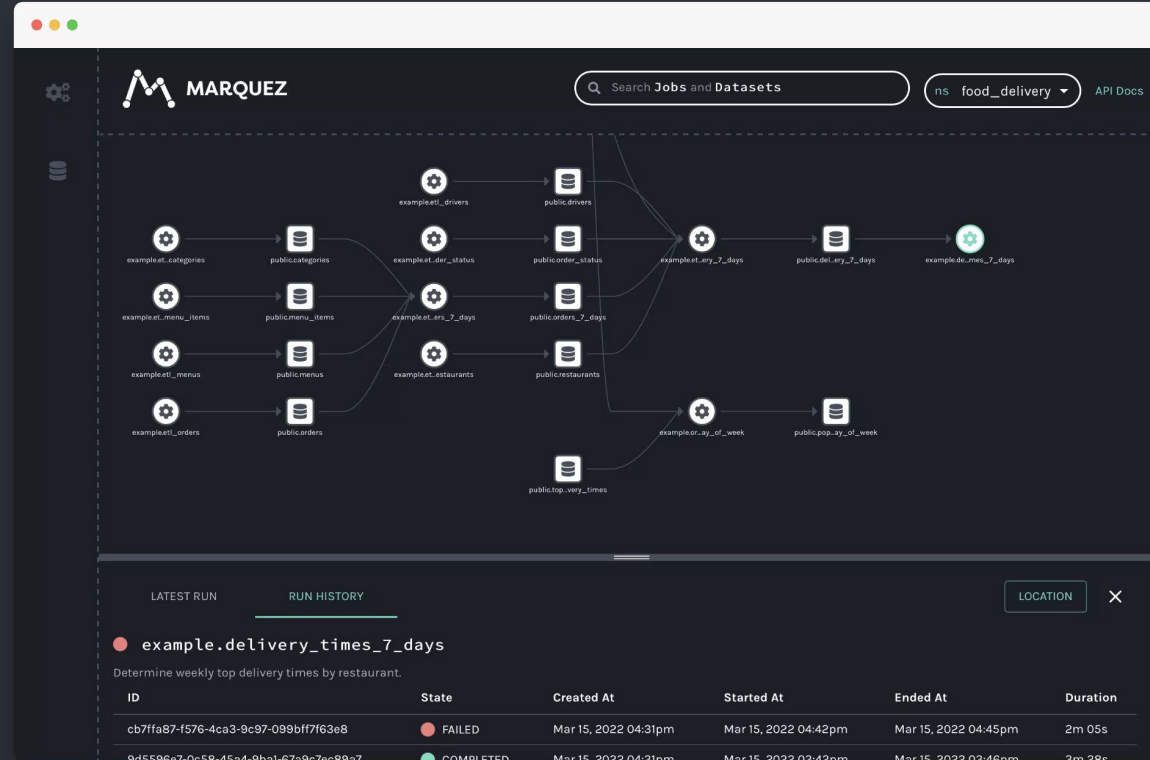
}



The screenshot shows the Apache Airflow web interface. At the top, there's a navigation bar with the Airflow logo, links for DAGs, Security, Browse, Admin, and Docs, and a clock showing 13:36 UTC. Below the navigation bar, the title 'DAGs' is displayed. A filter section shows 'All 30' DAGs, with 'Active 1' and 'Paused 29' counts. A search bar is labeled 'Search DAGs'. The main table lists DAGs with columns: DAG, Owner, Runs, Schedule, Last Run, Next Run, and Recent Tasks. The first DAG, 'example\_bash\_operator', is highlighted. Annotations with arrows point to specific elements: 'Names of workflows' points to the DAG name and its tags; 'Workflow schedules' points to the 'Schedule' column; 'State of workflow tasks' points to the 'Recent Tasks' column, specifically to the colored circles representing task states.

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
<input checked="" type="checkbox"/> example_bash_operator example example2	airflow	<input type="radio"/> <input checked="" type="radio"/> 2 <input type="radio"/>	0 0 ***	2021-11-06, 13:35:45	2021-11-06, 00:00:00	<input checked="" type="radio"/> 1 <input type="radio"/> 1 <input checked="" type="radio"/> 9 <input type="radio"/>
<input type="checkbox"/> example_branch_datetime_operator_2 example	airflow	<input type="radio"/> <input type="radio"/> <input type="radio"/>	@daily		2021-11-05, 00:00:00	<input type="radio"/> <input type="radio"/> <input type="radio"/>
<input type="checkbox"/> example_branch_dop_operator_v3 example	airflow	<input type="radio"/> <input type="radio"/> <input type="radio"/>	* * * *		2021-11-06, 13:34:00	<input type="radio"/> <input type="radio"/> <input type="radio"/>
<input type="checkbox"/> example_branch_labels	airflow	<input type="radio"/> <input type="radio"/> <input type="radio"/>	@daily		2021-11-05, 00:00:00	<input type="radio"/> <input type="radio"/> <input type="radio"/>
<input type="checkbox"/> example_branch_operator example example2	airflow	<input type="radio"/> <input type="radio"/> <input type="radio"/>	@daily		2021-11-05, 00:00:00	<input type="radio"/> <input type="radio"/> <input type="radio"/>
<input type="checkbox"/> example_complex example example2 example3	airflow	<input type="radio"/> <input type="radio"/> <input type="radio"/>	None			<input type="radio"/> <input type="radio"/> <input type="radio"/>
<input type="checkbox"/> example_dag_decorator example	airflow	<input type="radio"/> <input type="radio"/> <input type="radio"/>	None			<input type="radio"/> <input type="radio"/> <input type="radio"/>

# dataLineage.json





# Planejamento {

< Na etapa de **planejamento**, o objetivo é garantir que todas as partes envolvidas no ciclo de vida dos dados estejam alinhadas e que as necessidades do negócio sejam compreendidas antes de iniciar a implementação de pipelines ou sistemas. >

}

# Levantar o Objetivo {

- 1. Qual o objetivo de **negócio** do cliente?  
Dica: Importante o olhar para o negócio e não para o produto final.

Ex: Reduzir o tempo de processamento de pedidos em 30% nos próximos 12 meses, aumentando a satisfação do cliente e otimizando a eficiência operacional.

}

# Levantar Indicadores {

- 2. Quais indicadores podem medir este objetivo?

Ex:

- Tempo Médio de Processamento de Pedidos
- Net Promoter Score (Satisfação dos Clientes)

}

# Levantar Indicadores {

- 2.1. Quais fatores influenciam para que este objetivo seja atingido?

Ex:

- |                          |                          |
|--------------------------|--------------------------|
| - Capacidade             | - Automação de processos |
| computacional            | - Erro em transações     |
| - Disponibilidade de     | - Aprovação de Pagamento |
| Produtos                 | - Disponibilidade de     |
| - Nível de Demanda       | Sistemas                 |
| - Disponibilidade de mão | - Conformidade de        |
| de obra                  | Fornecedores             |

}

# Levantar Indicadores {

- 3. Crie métricas a partir dos indicadores propostos.

Ex:

- Tempo Médio de Processamento de Pedidos (Média)
- Quantidade de Pedidos (Soma)
- Percentual de uso dos recursos computacionais.
- Percentual de produtos disponíveis em estoque no momento em que o pedido é realizado.
- Percentual de fornecedores que entregam produtos dentro dos prazos e padrões acordados.

}

# Estudo de Fontes de Dados {

- 4. Qual a fonte de dados disponível para medir este indicador.

Ex:

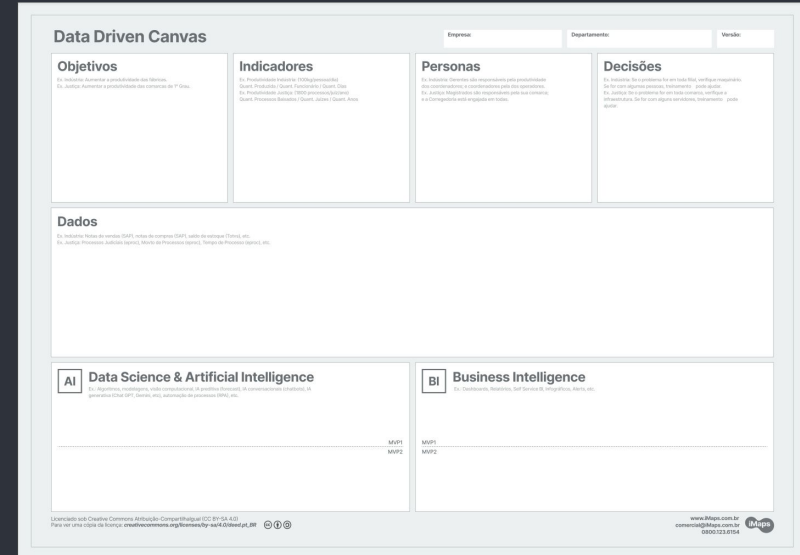
- Banco de Dados
- Planilha Excel
- Log de Sistema

}

# Documentação de Requisitos {

## Data Driven Canvas

< O DDC é uma abordagem visual e bastante rápida de documentar requisitos. >



# Documentação de Requisitos {

BUS Matrix

< O BUS Matrix ou Matriz Indicado X Dimensão é uma abordagem que confronta indicadores e dimensões em uma tabela >

}



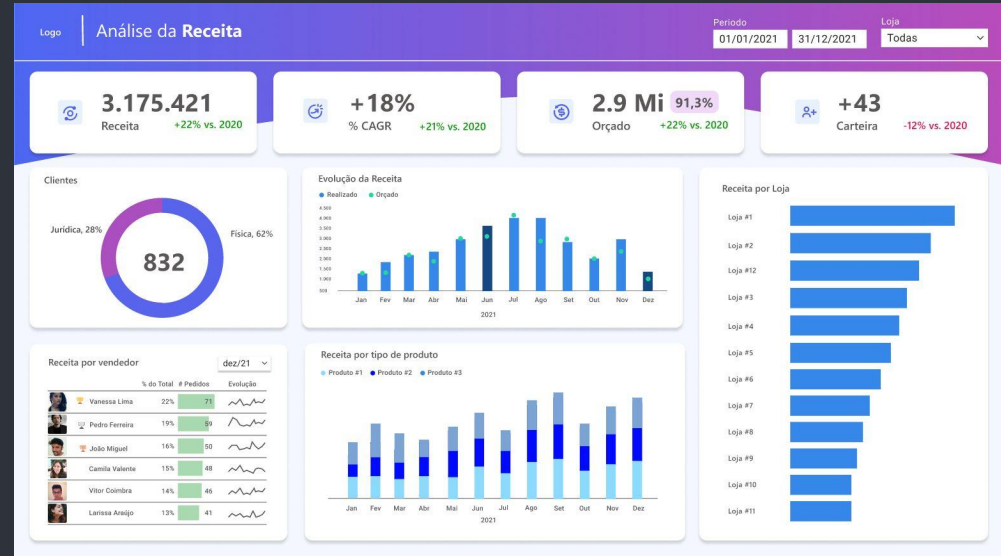
# Documentação de Requisitos {

## BUS Matrix

BUSINESS PROCESSES	SHARED DIMENSIONS									
	Account	Customer	Data	Department	Employee	Organization	Products	Promotion	Reseller	Sales Territory
Customer Service Calls		✓	✓		✓		✓			
Customer Surveys		✓	✓				✓			
General Ledger	✓		✓	✓		✓				
<b>Internet Sales</b>		✓	✓				✓	✓		✓
Inventory			✓				✓			
<b>Reseller Sales</b>			✓		✓		✓	✓	✓	✓
<b>Sales Plan</b>			✓				✓			✓

# Documentação de Requisitos {

Prototipação  
<Figma, PowerBi>



```
1
2
3 Business {
4
5     [Intelligence]
6
7
8
9
10
11
12
13
14
```

Prof. Juscelino Fernandes da Costa Junior

}

# Fontes de Dados {

## Estruturados



SGDB



API

## Não Estruturados



Planilha



Logs



Gmail



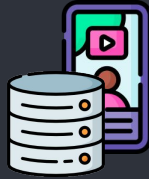
Vídeos



Fotos



Audio



## OLTP (Transacional)

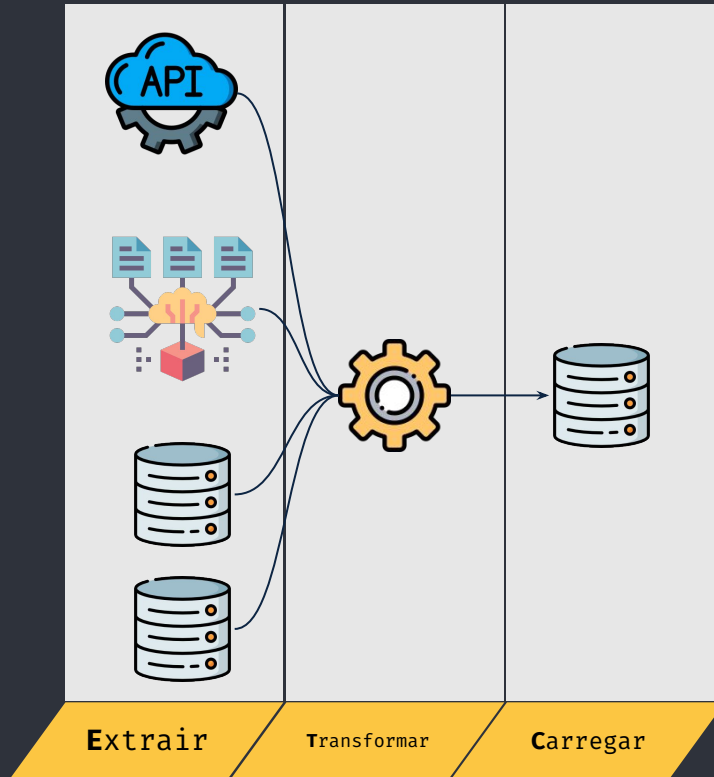
OLTP é otimizado para inserção, atualização e consistência dos dados em tempo real.



## OLAP (Analítico)

O OLAP é otimizado para leitura pesada e análises de dados complexos.

< Dashboards de B.I. devem evitar consumir dados diretamente de bancos OLTP >

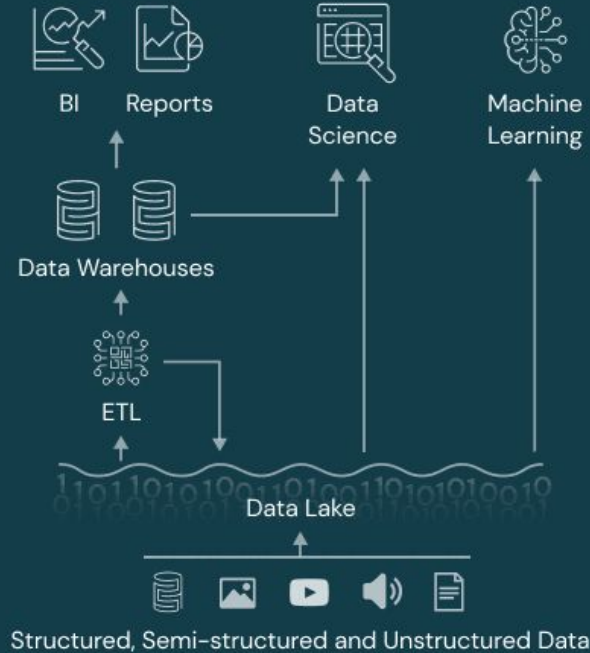


< **ETL** é o processo de Extrair, Transformar e Carregar dados. O intuito deste processo é extrair os dados de diversas fontes, aplicar as transformações necessárias e realizar cargas em um banco de dados OLAP, ao qual denominamos **Data Warehouse** >

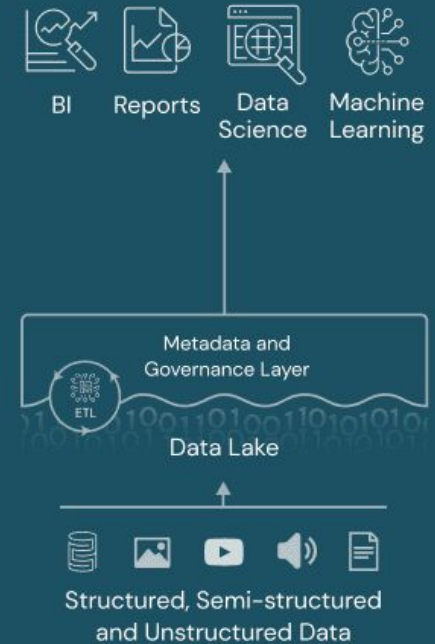
## Data Warehouse



## Data Lake



## Data Lakehouse



# Data Warehouse {

## Análises Multidimensionais

O DW organiza suas tabelas em **fatos** e **dimensões**.



### Fato

Uma tabela de fatos guarda **eventos numéricos** como vendas, lucros ou **quantidades**. Ela responde a "o quê, quanto e quando". É a base para análises e sempre se conecta a dimensões para dar contexto.

}



# Data Warehouse {

## Análises Multidimensionais

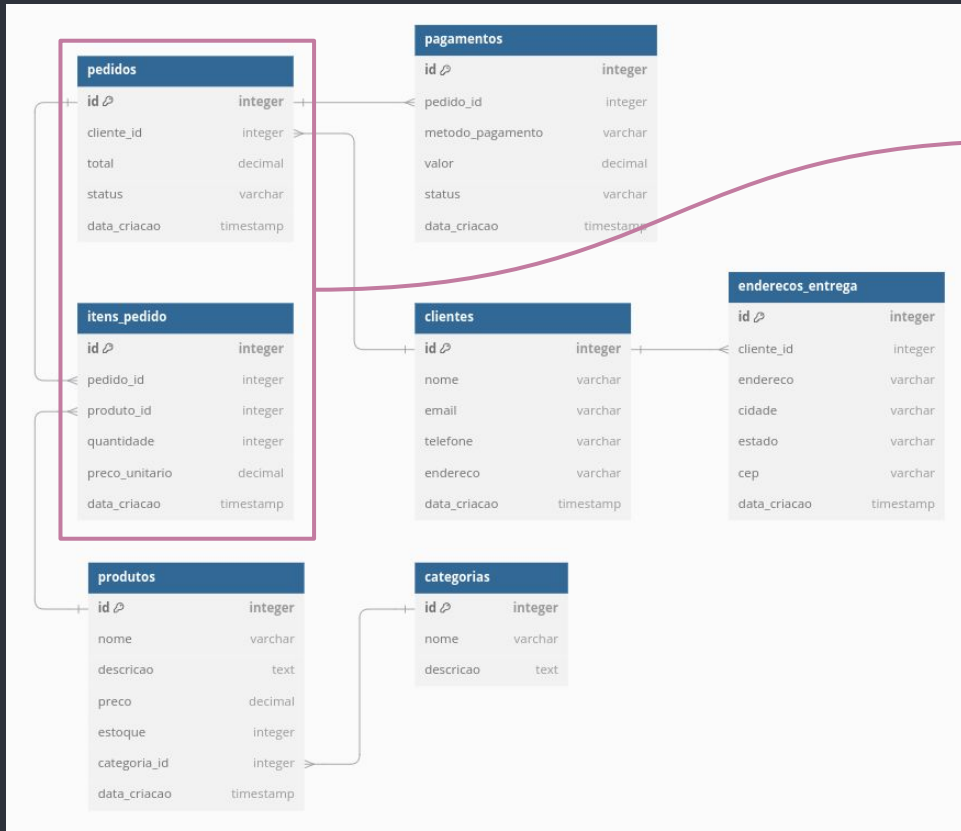
O DW organiza suas tabelas em **fatos** e **dimensões**.



## Dimensão

Uma tabela de dimensão organiza categorias, como tempo, produto ou cliente. Ela dá detalhes sobre o fato e permite explorar os dados por diferentes ângulos, como "quem, o quê e onde".

}

**Fatos:**

fVendas

fPagamentos

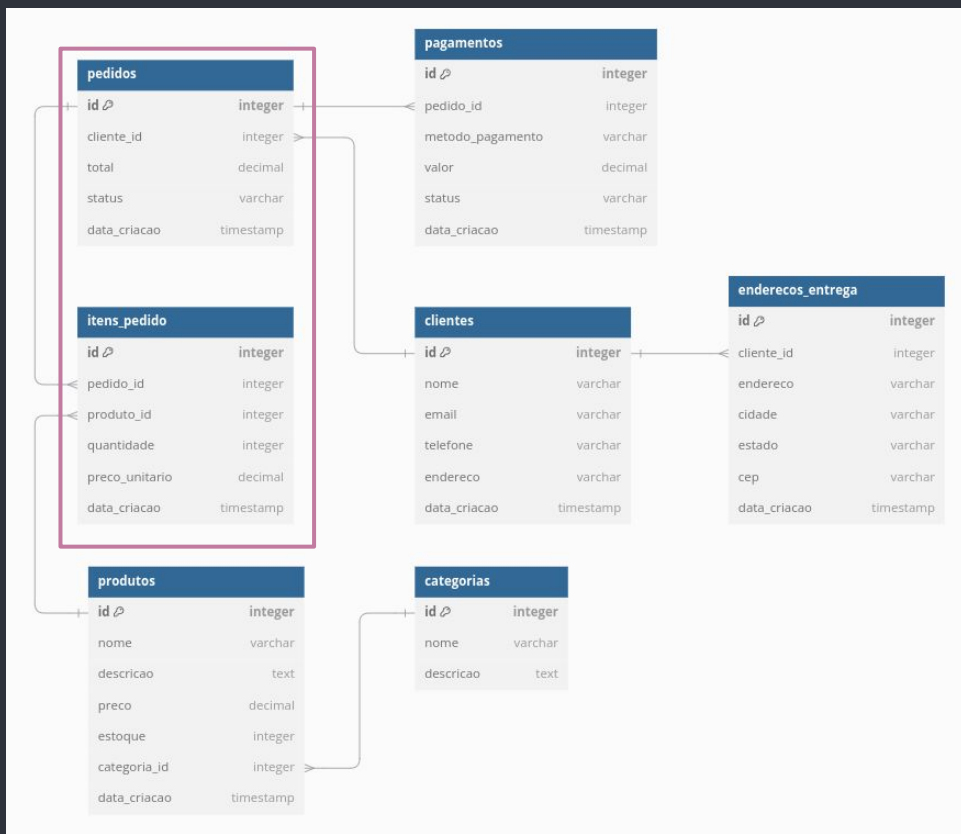
**Dimensões:**

dProduto

dCategoria

dCliente

dLocalidade



**Filtragens:**

status = 'VENDIDO'

**Cálculos e  
Agregações:**

$Vl\_venda =$   
 $quantidade *$   
 $preco\_unitario$

# Transformação {

## **Limpeza e Tratamento**

Se nulo → N/A , Remover linhas sem valores, correção ortográfica.

## **Agregação de Dados**

Soma, Média, Mediana, Desvio Padrão, Máximo, Mínimo, Concatenação.

## **Classificação**

Status = 1 → Pendente

Status = 2 → Concluído

## **Cálculos**

}

# BUS Matrix {

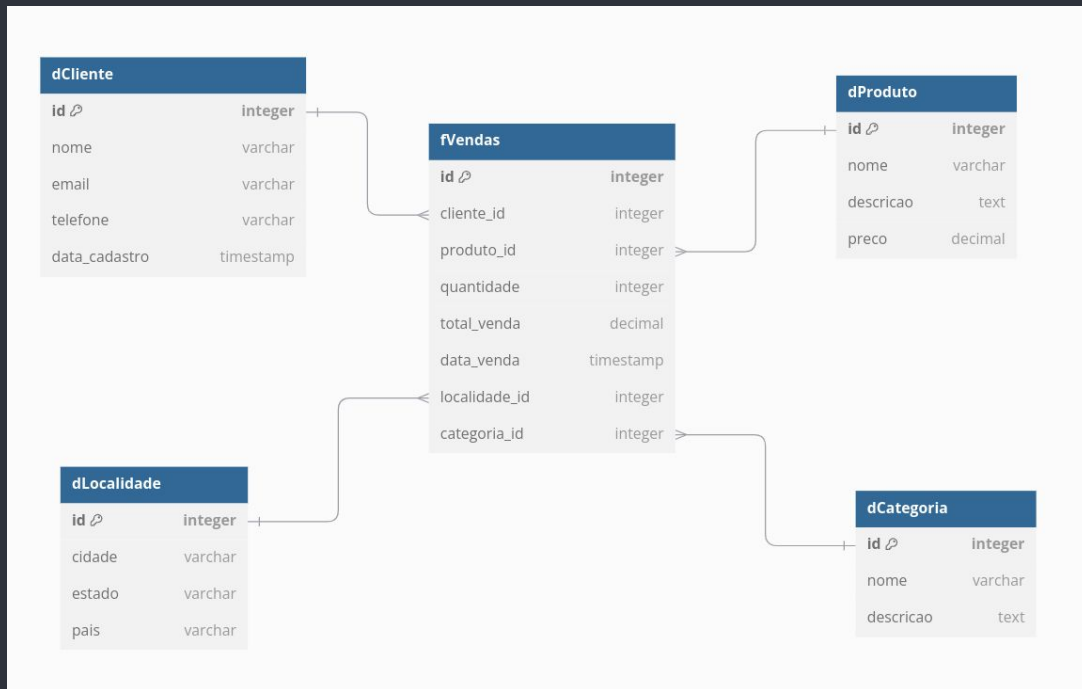
< O BUS Matrix,  
é uma matriz que  
confronta  
**indicadores** e  
**dimensões**,  
facilita  
bastante para  
tratar  
diretamente com  
o cliente. >

Indicador	Dimensões				
	Produto	Vendedor	Fornecedor	Cliente	Calendario
Qtd Vendas	x	x		x	x
Qtd Vendedores		x			
Custo Matéria	x		x		
Qtd Reclamacoes	x	x		x	x

# Star Schema {

< O Star Schema ou Estrela, é um formato de modelagem de dados mais utilizado em Data Warehouses.  
>

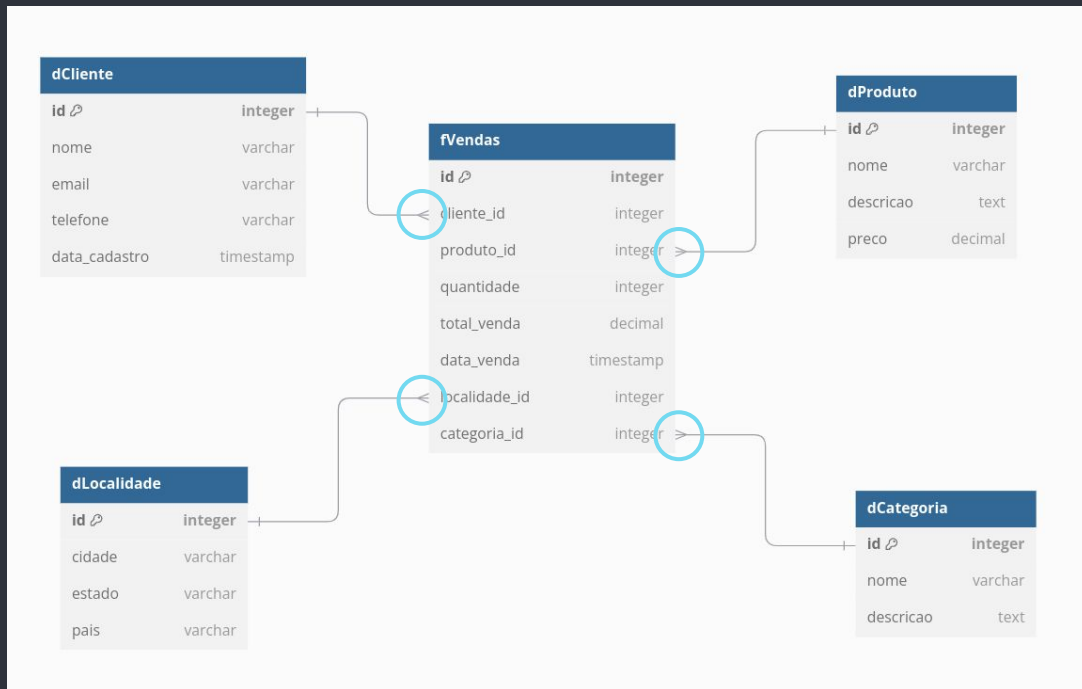
}



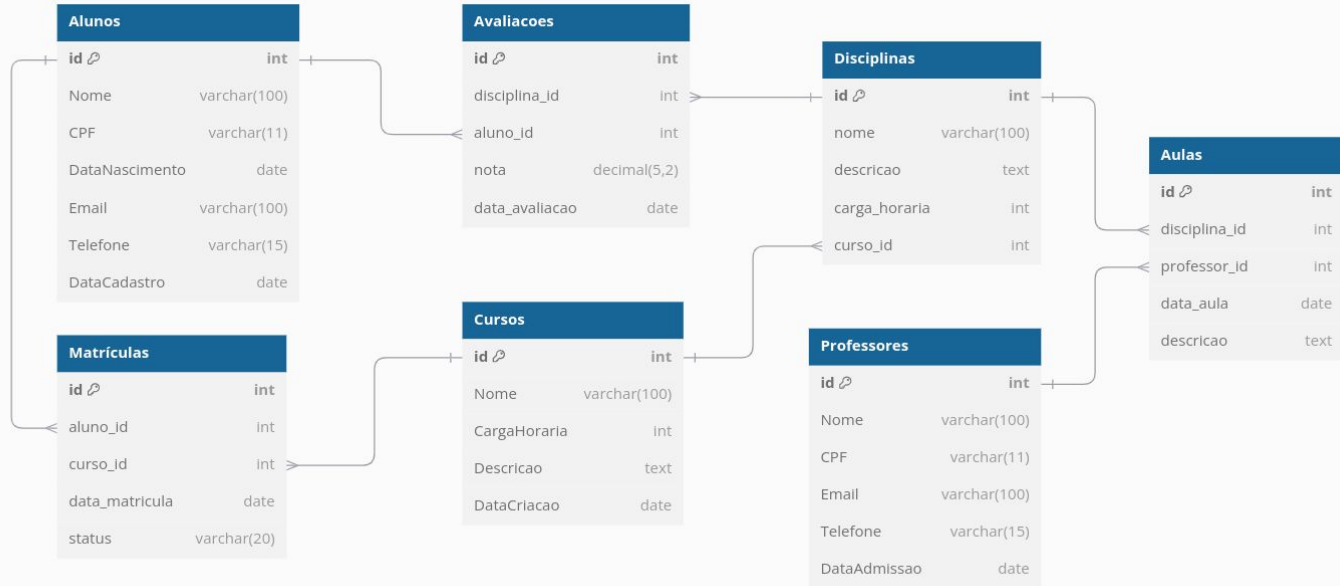
# Star Schema {

< Uma característica marcante no Star Schema, é o relacionamento n:1 da fato para as dimensões. Um modelo estrela deve sempre seguir este mesmo formato. >

}



## TDE {





# Business { [Intelligence]

Prof. Juscelino Fernandes da Costa Junior

}

# ETL < Extração />



/\*\* A **extração** é a base de um pipeline ETL, é o processo de leitura dos dados. \*/

/\*\* Vamos falar de pipeline mais adiante. \*/



## Batch

Processamento em lote (batch) é a execução de tarefas ou cargas de trabalho de forma agrupada e programada, processando grandes volumes de dados de uma vez.

**Utilização:** Ideal para processamento periódico de grandes quantidades de dados, onde o tempo de resposta imediato não é essencial.

### Principais Ferramentas:

- Injestores:
  - Spark
  - Airbyte
- Orquestradores:
  - Airflow
  - Jenkins



## Streaming

Processamento em streaming é o tratamento de fluxos contínuos de dados em **tempo real**. Os dados são processados à medida que chegam, permitindo respostas e ações rápidas baseadas em eventos.

**Utilização:** Ideal para cenários onde os dados precisam ser processados continuamente.

### Principais Ferramentas:

- Kafka
- Spark
- Dataflow (gcp)



## Full Extraction

Extrai **todos** os dados da fonte.

Útil em situações onde os dados não mudam frequentemente, mas exige mais tempo e processamento.



## Incremental

**Apenas** os dados que foram **alterados** ou **adicionados** desde a última extração são extraídos.

Reduz o volume de dados e otimiza o processo de ETL.

1

2

```
# Consulta para extrair todos os dados da tabela de vendas
query_full = "SELECT * FROM sales"
df_sales_full = pd.read_sql(query_full, mysql_engine)

# Carregando os dados para o Snowflake
df_sales_full.to_sql(
    'sales_snowflake',
    snowflake_engine,
    index=False,
    if_exists='replace'
)
```

Python

13

14

```
# Definindo a data da última carga incremental
ultima_data_carga = '2024-10-10 12:00:00'

# Consulta para extrair apenas dados novos ou modificados após a última carga
query_incremental = f'SELECT * FROM sales WHERE last_modified > {ultima_data_carga}'

df_sales_incremental = pd.read_sql(query_incremental, mysql_engine)

# Carregando os dados incrementais para o Snowflake
df_sales_incremental.to_sql(
    'sales_snowflake',
    con=snowflake_engine,
    index=False,
    if_exists='append'
)
```