#### Análise de Dados com Base em Processamento Massivo em Paralelo

# Aula 4: Modelagem Conceitual de ETL/ELT

Cristina Dutra de Aguiar Ciferri ICMC/USP cdac@icmc.usp.br









Características

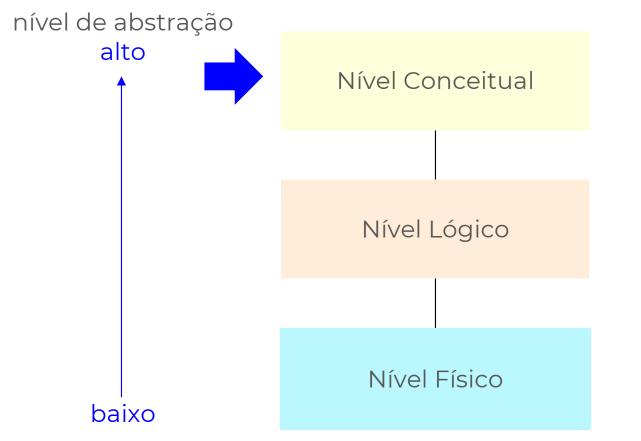
Modelo Intuitive

• Exemplo para a BI Solutions

#### Projeto de ETL/ELT

- Características desejadas para o sucesso
  - Robustez
  - Boa documentação
  - Facilidade de Manutenção
- Representado como um workflow
  - Cadeia de operações ou tarefas aplicadas aos dados
  - o Representado por meio de um modelo

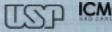
### Níveis de Abstração de um Modelo



complementa a análise de requisitos, facilitando o entendimento do processo

descreve os detalhes técnicos das tarefas envolvidas

incorpora aspectos de implementação e otimização







#### **Desafios e Motivação**

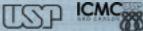
- Desafios
  - Criticidade e complexidade do processo de ETL/ELT
  - Grande esforço despendido para a construção do processo
  - Propensão a falhas
- Motivação
  - o Facilitar e padronizar a construção do processo de ETL/ELT
  - Melhorar a qualidade do processo de ETL/ELT e dos dados armazenados no DW

#### **Modelagem Conceitual**

- Realizada na fase inicial do processo de ETL/ELT
  - Requisitos dos usuários de SSD
  - Entendimento do conteúdo e da estrutura das fontes de dados
  - Enfoque na estrutura proposta para o DW
- Produz um esquema conceitual
  - Representação gráfica e abstrata do processo de ETL/ELT

#### Requisitos da Modelagem Conceitual

- Características desejadas
  - Simplicidade e completude
  - Clareza, consistência, não ambiguidade
- Diagrama produzido
  - Deve ser facilmente entendido pelos usuários finais que são conhecedores do negócio
    - muitas vezes não possuem conhecimento profundo de tecnologias
  - o Deve contribuir para diminuir o esforço dos projetistas e desenvolvedores





#### **Funcionalidades Adicionais**

- Documenta as decisões tomadas
- Possibilita a análise de impacto das alterações que ocorrem no ciclo de vida da aplicação de data warehousing
  - Alterações nas fontes de dados
  - Evolução dos requisitos ou das regras de negócio
  - Correção de erros cometidos durante a fase de projeto
- Facilita a exploração de cenários alternativos



#### Modelagem Conceitual versus Ferramentas

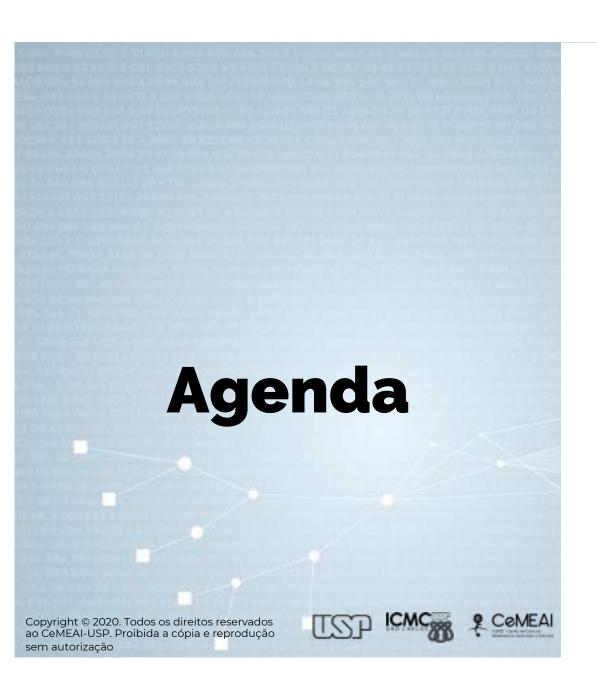
- Modelagem Conceitual
  - o Fornece alto nível de abstração, sendo independente de ferramentas específicas
- Ferramentas de ETL/ELT disponíveis
  - Relacionadas aos níveis lógico e físico
  - Exemplos
    - Pentaho Data Integration (Kettle)
    - Talend Open Studio
    - CloverETL

- Oracle Warehouse Builder
- IBM Infosphere
- MSSQLServer Integration Services









Características

Modelo Intuitive

• Exemplo para a BI Solutions

#### Características do Modelo

- Operadores
  - Entidades de alto nível
  - Representam as operações típicas de ETL/ELT
  - Possuem notação gráfica
- Combinação de operadores
  - o Realizada por setas unidirecionais que indicam a propagação dos dados
  - o Representa sequências que compõem o workflow de ETL/ELT







#### Características do Workflow

- Início e final
  - Início: um ou mais repositórios de dados
  - o Final: um ou mais repositórios, sendo o principal o data warehouse (ou data mart)

- Funcionalidades dos operadores
  - Manipulação de dados
  - Organização do fluxo de dados no workflow

Entrada Parâmetro Saída





### Especificação dos Operadores: Entrada

Um ou mais conjuntos de dados

- Classificação
  - Unária: apenas um conjunto de dados
  - o Binária: dois conjuntos de dados
  - N-ária: vários conjuntos de dados





# Especificação dos Operadores: Tipos de Parâmetro

- Lista de atributos
  - Nome de um atributo do conjunto de dados
  - o Exemplo: funcNome, funcMatricula
- Condição
  - Expressão relacional
    - Exemplo: funcCidade = São Carlos
  - Expressão lógica
    - funcEstadoSigla = SP AND funcMatricula > 32879

Operações relacionais



Operadores lógicos NOT, AND, OR







# Especificação dos Operadores: Tipos de Parâmetro

- Ordenação
  - Ordem crescente ou decrescente dos dados
  - Exemplo: asc e desc
- Precedência
  - Qual conjunto de dados deve ser analisado primeiro
  - Exemplo: dados do conjunto A dados do conjunto B
- Lista de atribuições
  - Atributo <--- valor
  - Exemplo: funcMatricula <--- 234334, funcEstadoSigla <--- PE

#### Especificação dos Operadores: Saída

Um ou mais conjuntos de dados

- Classificação
  - Unária: apenas um conjunto de dados -
  - o Binária: dois conjuntos de dados
  - N-ária: vários conjuntos de dados





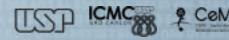


#### Categorias de Operadores

- Classificação baseada em
  - Características dos operadores
  - o Efeitos que causam sobre os dados ou sobre a organização do processo

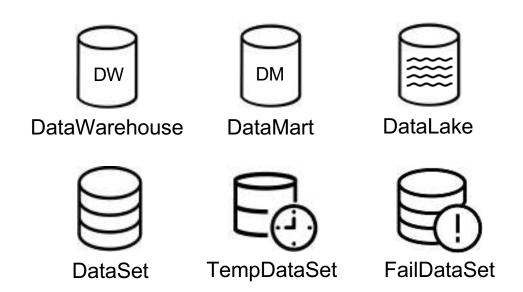
Armazenamento Manipulação de Dados Inicialização

Agregação Fluxo Especiais



#### Operadores de Armazenamento

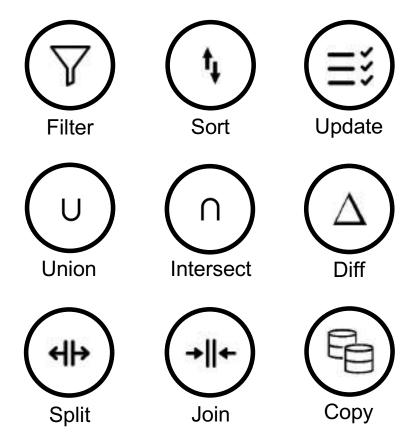
Representam locais de armazenamento de dados, tais como repositórios, arquivos ou bancos de dados







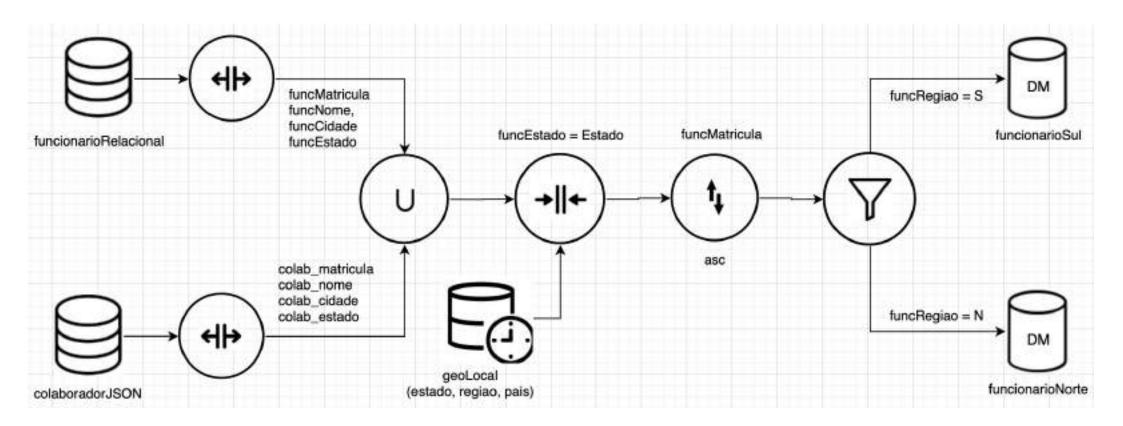
Representam operações de transformação e de limpeza dos dados







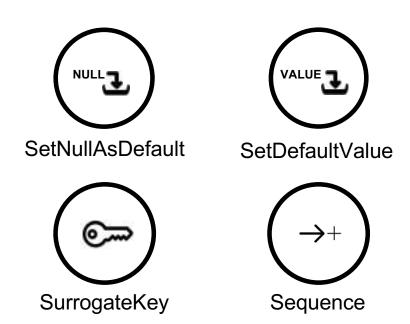
#### Geração de Data Marts Regionais de Funcionários





#### Operadores de Inicialização

Representam a inicialização de um atributo com um valor específico

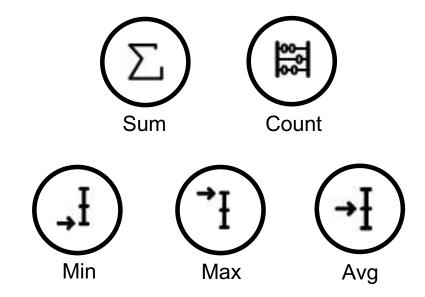






### Operadores de Agregação (1/2)

Representam funções que processam os valores de um atributo e retornam um único valor







### Operadores de Agregação (2/2)

Representam funções que processam os valores de um atributo e retornam um único valor para cada atributo do agrupamento









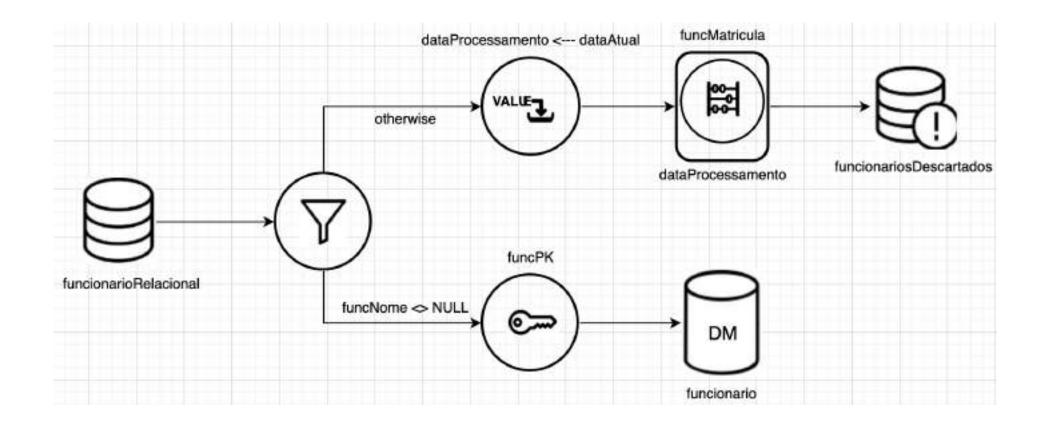








#### Análise do Processamento de Funcionários

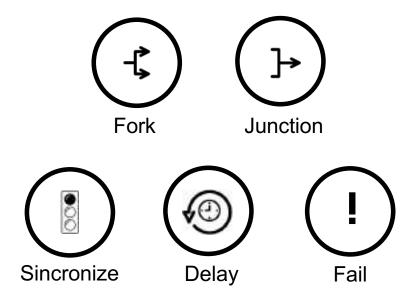






### Operadores de Fluxo

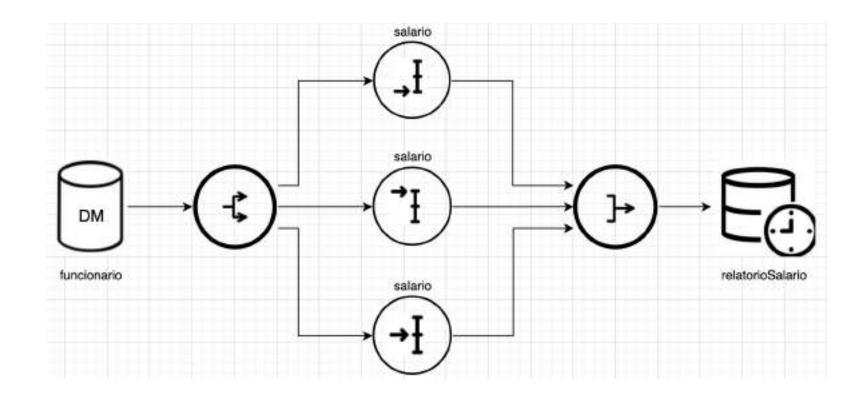
Representam uma alteração no fluxo dos dados, sem impactar esses dados







## Geração de Relatório de Salários



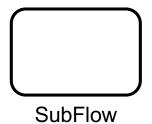


#### **Operadores Especiais**

Representam operações
que envolvem
especificidades,
complementando as
funcionalidades dos demais
operadores



Function title
Short description
Details









#### **Material Suplementar**

- Tabela descritiva dos operadores
  - Operador e sua funcionalidade
  - Representação gráfica, incluindo entrada, parâmetro e saída

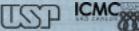
Operador	Funcionalidade	Representação Gráfica
Fork	direcionar um conjunto de dados para duas ou mais tarefas paralelas	···O—••••••••••••••••••••••••••••••••••

Acompanha os slides no formato .pdf



#### Operadores de Armazenamento

Operador	Funcionalidade	Representação Gráfica
DataWarehouse	armazena dados multidimensionais	DW ou DW nome
DataMart	data warehouse com escopo limitado	DM OU DM nome
DataLake	armazena dados brutos que ainda não foram transformados	→ a ou a nome nome







#### Operadores de Armazenamento

Operador	Funcionalidade	Representação Gráfica
DataSet	armazena dados	ou on nome
TempDataSet	área temporária de armazenamento de dados	ou nome
FailDataSet	armazena dados rejeitados por uma operação ou para efeitos de log	ou nome







Operador	Funcionalidade	Representação Gráfica
Filter	seleciona subconjuntos de dados de acordo com condições definidas	condiçãol  condiçãoN otherwise
Sort	ordena dados em ordem crescente ou decrescente, de acordo com atributos definidos	atributol asc,, atributoN desc
Update	altera os valores dos dados, de acordo com condições definidas sobre atributos	lista de atribuições  condição



Operador	Funcionalidade	Representação Gráfica
Union	une conjuntos de dados, gerando um conjunto que contém todos os dados de entrada, sem repetição	
Intersect	une conjuntos de dados, gerando um conjunto que contém apenas os dados em comum, sem repetição	O otherwise
Diff	gera os dados que estão presentes no primeiro conjunto de dados, mas não estão no segundo conjunto	$\begin{array}{c} conjunto1 \\ \hline \\ conjunto2 \\ \hline \\ conjunto1 - conjunto2 \\ \end{array}$



Operador	Funcionalidade	Representação Gráfica
Split	separa atributos de um conjunto de dados, direcionando-os para fluxos diferentes	atributol, atributo2 atributoN
Join	combina dois conjuntos de dados usando como base atributos em comum	conjuntol.atributol = conjunto2.atributo2
Сору	a partir de um conjunto de dados de entrada, gera o próprio conjunto e uma réplica deste	A réplica de A



### Operadores de Inicialização

Operador	Funcionalidade	Representação Gráfica
SetNullAsDefault	inicializa um atributo específico com o valor nulo, para todos os itens do conjunto de dados	lista de atributos  NULL  condição
SetDefaultValue	atribui um determinado valor para um atributo específico, para todos os todos os itens do conjunto de dados	lista de atribuições

### Operadores de Inicialização

Operador	Funcionalidade	Representação Gráfica
SurrogateKey	cria um atributo chave e atribui a ele um valor único para cada item do conjunto de dados	atributo
Sequence	cria um atributo não-chave e atribui a ele um valor único para cada item do conjunto de dados, o qual é gerado a partir de um valor inicial	atribuição →+

### Operadores de Agregação

Operador	Funcionalidade	Representação Gráfica
Sum	para cada atributo, soma seus valores e produz um único valor	lista de atributos
Count	para cada atributo, conta seus valores e produz um único valor	lista de atributos  condição

# Operadores de Agregação

Operador	Funcionalidade	Representação Gráfica
Min	para cada atributo, produz o menor valor	lista de atributos  condição
Max	para cada atributo, produz o maior valor	lista de atributos  condição
Avg	para cada atributo, produz o valor médio	lista de atributos  condição



#### Operadores de Agregação

Operador	Funcionalidade	Representação Gráfica
SumGroup	para cada grupo, soma os valores dos dados de cada atributo	lista de atributos  lista de atributos de agrupamento condição
Count	para cada agrupamento, conta o número de dados de cada atributo	lista de atributos  lista de atributos de agrupamento condição

# Operadores de Agregação

Operador	Funcionalidade	Representação Gráfica
MinGroup	para cada grupo, produz o menor valor de cada atributo	lista de atributos  Lista de atributos de agrupamento condição
MaxGroup	para cada grupo, produz o maior valor de cada atributo	lista de atributos  Iista de atributos  lista de atributos de agrupamento condição
AvgGroup	para cada grupo, produz o valor médio de cada atributo	lista de atributos  lista de atributos de agrupamento condição



# Operadores de Fluxo

Operador	Funcionalidade	Representação Gráfica
Fork	direciona um conjunto de dados para dois ou mais fluxos executados em paralelo ou para um repositório e fluxos	
Junction	junta dois ou mais fluxos executados em paralelo	

# Operadores de Fluxo

Operador	Funcionalidade	Representação Gráfica
Sincronize	sincroniza dois ou mais fluxos paralelos com base em uma condição de finalização	condição N
Delay	temporiza o tempo no qual será feita a análise de conjuntos de dados de fluxos paralelos	condiçãoN tempo
Fail	representa um fluxo alternativo para indicar falha	<u></u>







# **Operadores Especiais**

Operador	Funcionalidade	Representação Gráfica
Function	representa operações ou atividades muito específicas que não podem ser representadas pelos outros operadores	Function title Short description Details  nome
SubFlow	encapsula subfluxos que envolvem conjuntos de tarefas específicas	nome rótuloSaídaM rótuloSaídaM rótuloSaídaM rótuloSaídaM rótuloSaídaM



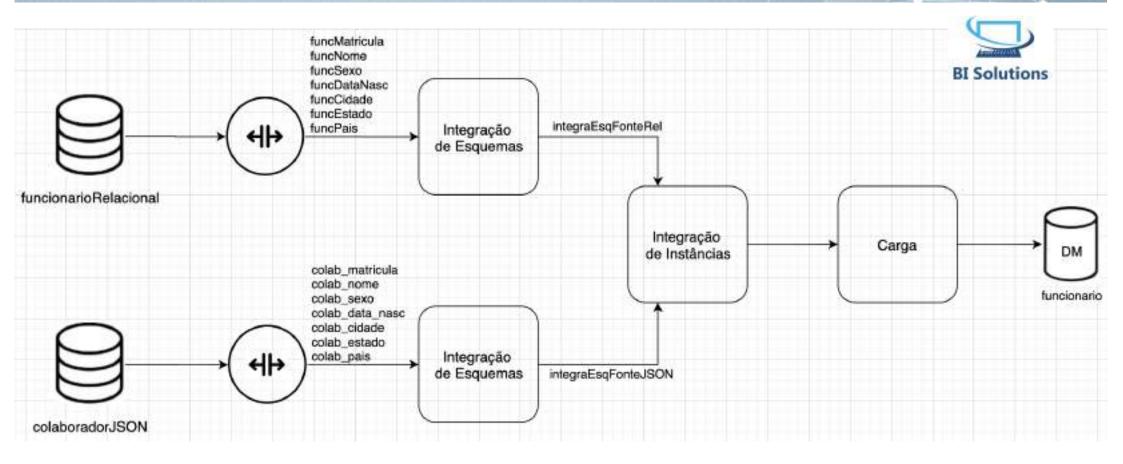


Características

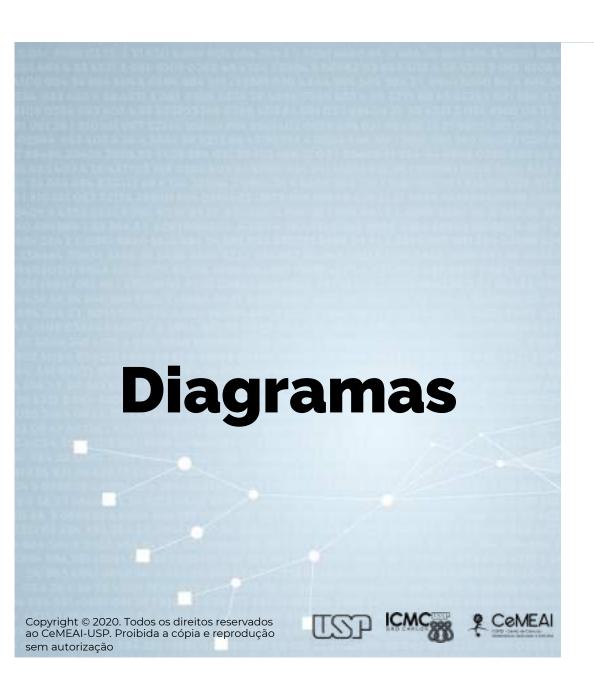
• Modelo Intuitive

• Exemplo para a BI Solutions

#### Processo de ETL da BI Solutions





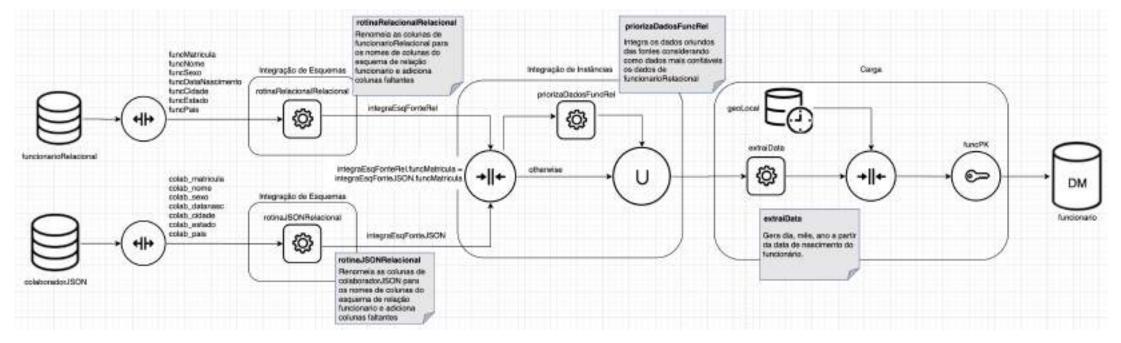


Exemplo do Processo de ETL

• Implementação em Pandas

#### **Diagrama Conceitual Completo**



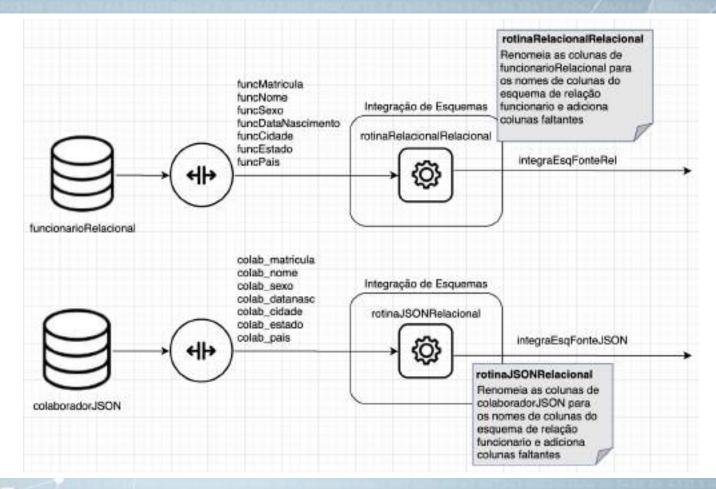








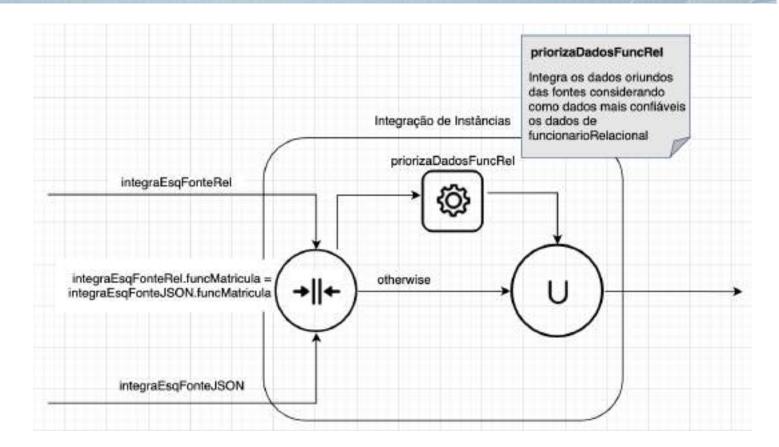
#### Extração e Integração de Esquemas







#### Integração de Instâncias

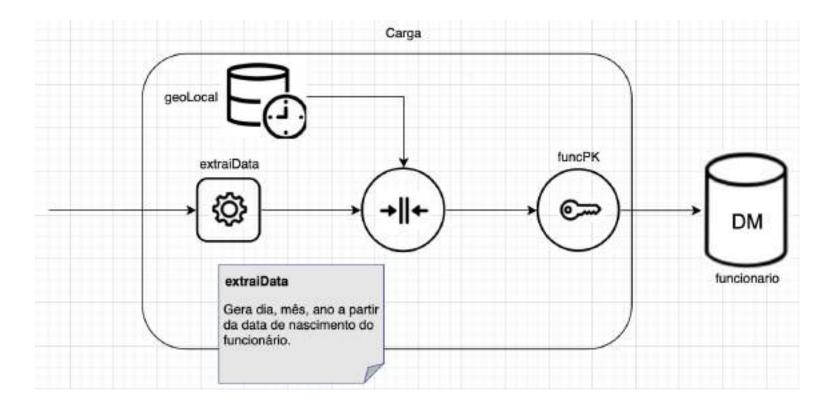






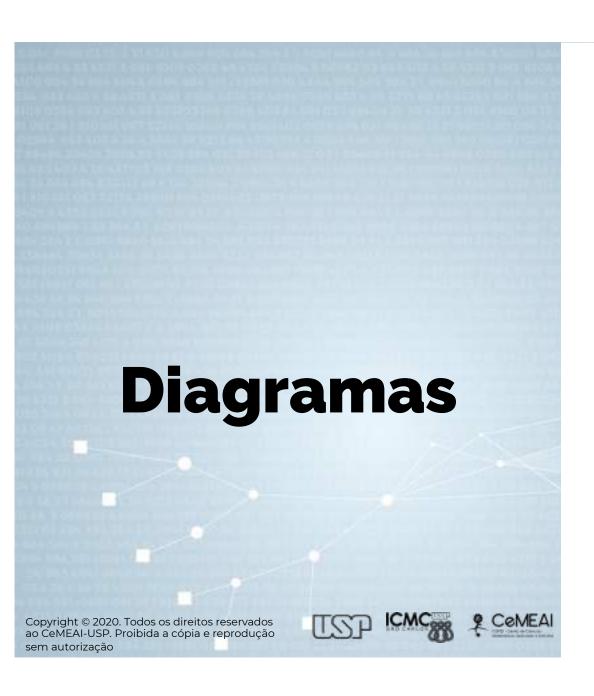
#### Carga









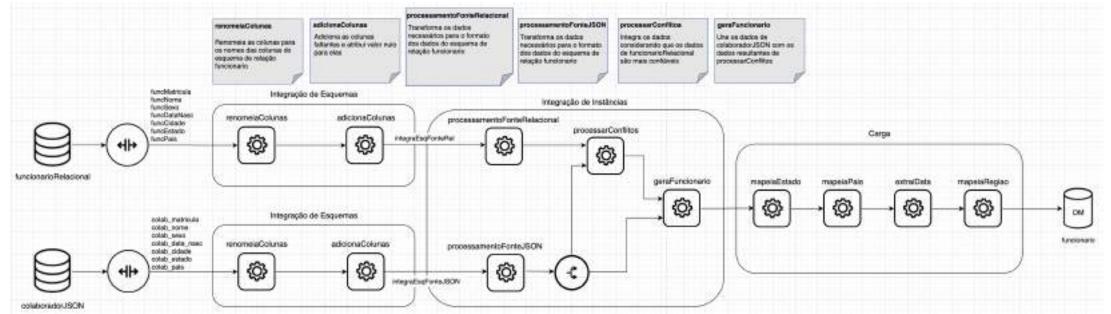


Exemplo do Processo de ETL

• Implementação em Pandas

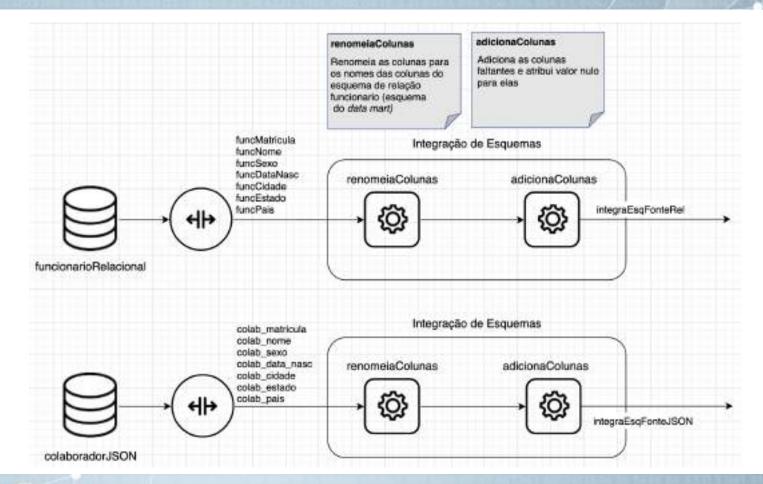
# Visão Geral da Implementação em Pandas







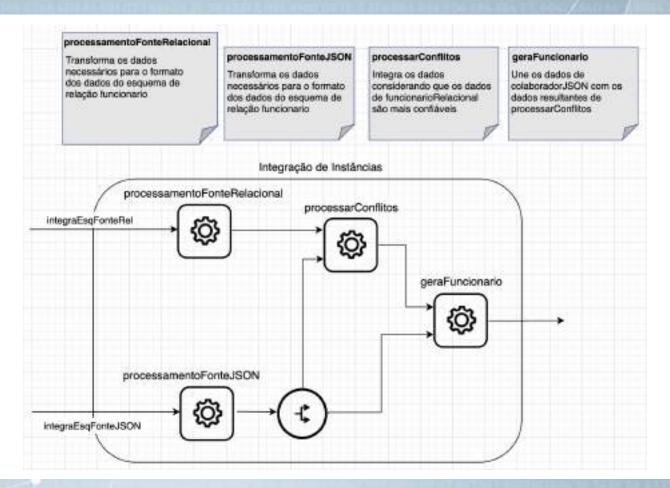
## Extração e Integração de Esquemas







# Integração de Instâncias









# Carga



