

# Unicamp - MC536 - Projeto Final

**Equipe Dinossauros, Bancos de dados e Coisas Parecidas (DBDCP)**

Davi Gabriel Bandeira Coutinho (183710)  
Francisco Vinicius Sousa Guedes (260440)  
Márcio Levi Sales Prado (183680)

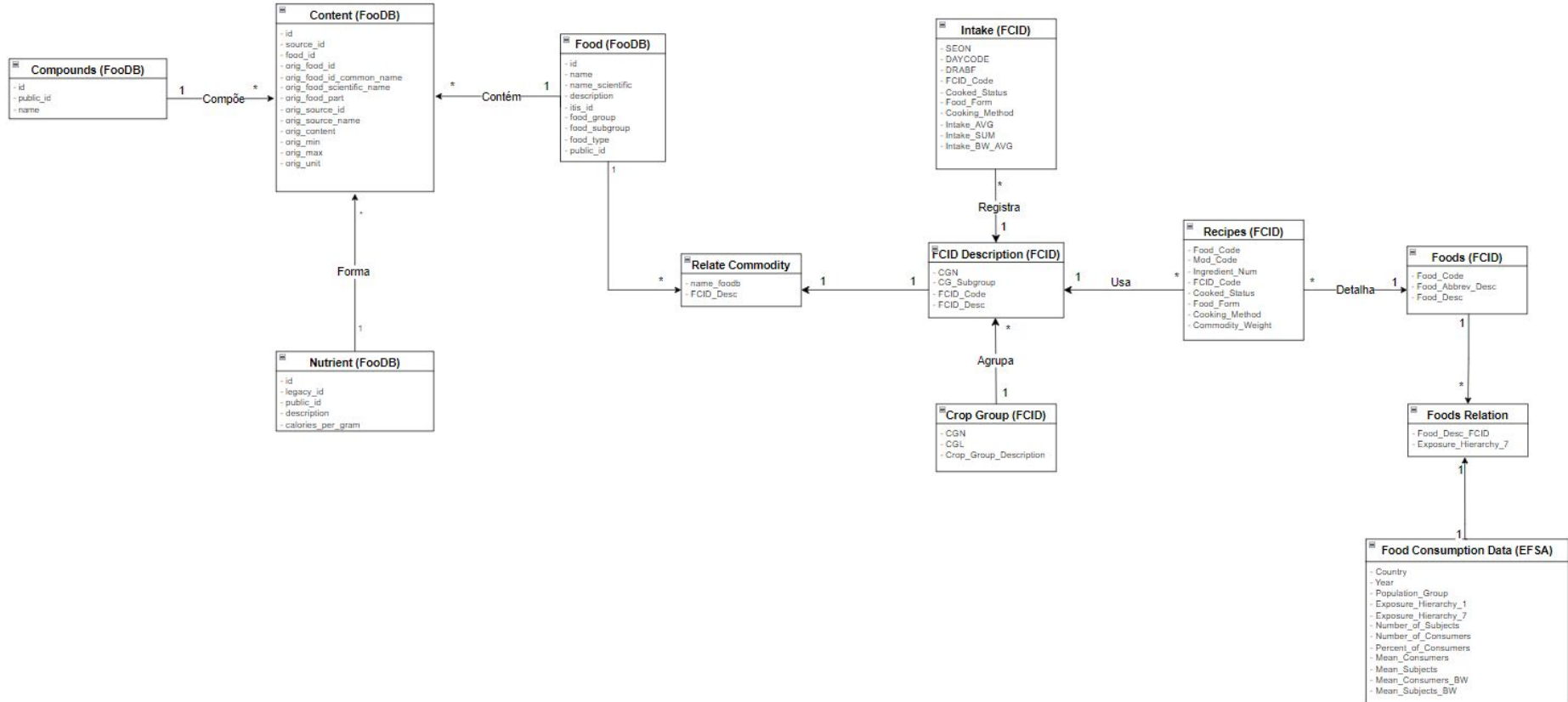
# Motivação e Objetivo

Com a constante aceleração do dia a dia, muitas pessoas optam por preencher suas refeições com comidas rápidas e processadas, fazendo **escolhas não muito saudáveis**.

O desenvolvimento e estudo dos **hábitos alimentares da população** também trazem à tona outras questões, desta vez relacionadas à saúde dos consumidores já que a alimentação de uma pessoa pode nos dizer bastante sobre sua saúde, como por exemplo indicar a **deficiência de algum nutriente**.

Nosso objetivo então é :  
**fornecer informações nutricionais sobre os hábitos alimentares da população alvo da pesquisa**, realizando uma conexão entre as bases de dados de consumo com o FooDB. Também buscaremos **trazer informações sobre receitas de alimentos**, a partir da análise dos dados que detalham os ingredientes que as compõem.

# Modelo Lógico Conceitual



# Modelo Lógico Relacional (Parte 1)

FOOD\_CONSUMPTION\_DATA(\_Country\_, \_Year\_, \_Population\_Group\_, Exposure\_Hierarchy\_1, \_Exposure\_Hierarchy\_7\_,  
Number\_of\_Subjects, Number\_of\_Consumers, Percent\_of\_Consumers, Mean\_Consumers, Mean\_Subjects, Mean\_Consumers\_BW,  
Mean\_Subjects\_BW)

FOODS\_RELATION(\_Food\_Desc\_FCID\_, \_Exposure\_Hierarchy\_7\_)

Food\_Desc\_FCID chave estrangeira -> FOODS(Food\_Desc)

Exposure\_Hierarchy\_7 chave estrangeira -> FOOD\_CONSUMPTION\_DATA(Exposure\_Hierarchy\_7)

FOODS(\_Food\_Code\_, Food\_Abbrev\_Desc, Food\_Desc)

RECIPES(\_Food\_Code\_, Mod\_Code, Ingredient\_Num, \_FCID\_Code\_, Cooked\_Status, Food\_Form, Cooking\_Method, Commodity\_Weight)

FCID\_DESCRIPTION(\_FCID\_Code\_, FCID\_Desc, CGN, CG\_Subgroup)

INTAKE(\_SEQN\_, \_DAYCODE\_, DRABF, FCID\_Code, Cooked\_Status, Food\_Form, Cooking\_Method, Intake\_AVG, Intake\_SUM,  
Intake\_BW\_AVG)

FCID\_Code chave estrangeira -> FCID\_DESCRIPTION(FCID\_Code)

# Modelo Lógico Relacional (Parte 2)

CROP\_GROUP(\_CGN\_, \_CGL\_, Crop\_Group\_Description)

CGN chave estrangeira -> FCID\_DESCRIPTION(CGN)

RELATE\_COMMODITY(\_name\_foodb\_, \_FCID\_Desc\_)

name\_foodb chave estrangeira -> FOOD(name)

FCID\_Desc chave estrangeira -> FCID\_DESCRIPTION(FCID\_Desc)

FOOD(\_id\_, name, name\_scientific, description, itis\_id, food\_group, food\_subgroup, food\_type, public\_id)

NUTRIENT(\_id\_, legacy\_id, public\_id, description, calories\_per\_gram)

COMPOUNDS(\_id\_, public\_id, name)

CONTENT(\_id\_, source\_id, \_food\_id\_, orig\_food\_id, orig\_food\_id\_common\_name, orig\_food\_id\_scientific\_name, orig\_content, orig\_min, orig\_max, orig\_unit)

id chave estrangeira -> NUTRIENT(id)

id chave estrangeira -> COMPOUNDS(id)

food\_id chave estrangeira -> FOOD(id)

# Matches entre Bases Diferentes (TF-IDF)

#1: Iniciamos carregando as 2 tabelas que vão ser comparadas.

#2: Desconsideramos letras maiúsculas, igualando-as a letra minúsculas

#3: Criamos o dicionário a partir das palavras que aparecem na base escolhida por referência

#4: Fazemos o encaixe das palavras no dicionário

#5: Criamos o vetor de presença de palavras

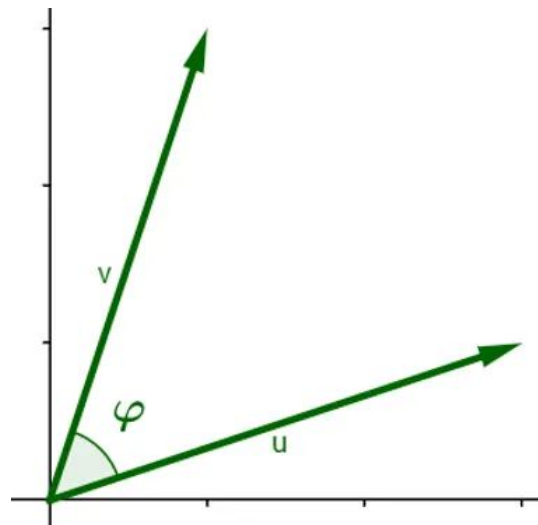
#6: Seleciona o melhor vetor, utilizando o `arg_max` sobre a similaridade do cosseno dada por `cosine_similarity(v1, v2)`

Count Vectorizer

	blue	bright	sky	sun
Doc1	1	0	1	0
Doc2	0	1	0	1

TD-IDF Vectorizer

	blue	bright	sky	sun
Doc1	0.707107	0.000000	0.707107	0.000000
Doc2	0.000000	0.707107	0.000000	0.707107



# Matches entre Bases Diferentes (TF-IDF)

Carregamento

Normalização

Dicionário e Conversão

Similaridade de  
Cosseno

Selecionando o  
melhor match

Gerando output

```
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

# Carregar dados do FCID e do FoodDB (Configurar o ambiente para estar na raiz do repositório.)
fcid_data = pd.read_csv('bases/fcid/FCID_Code_Description.csv')
fooddb_data = pd.read_csv('bases/fooddb/Food.csv')

# Pré-processamento: converter todos os nomes para minúsculas
fcid_data['processed'] = fcid_data['FCID_Desc'].str.lower()
fooddb_data['processed'] = fooddb_data['name'].str.lower()

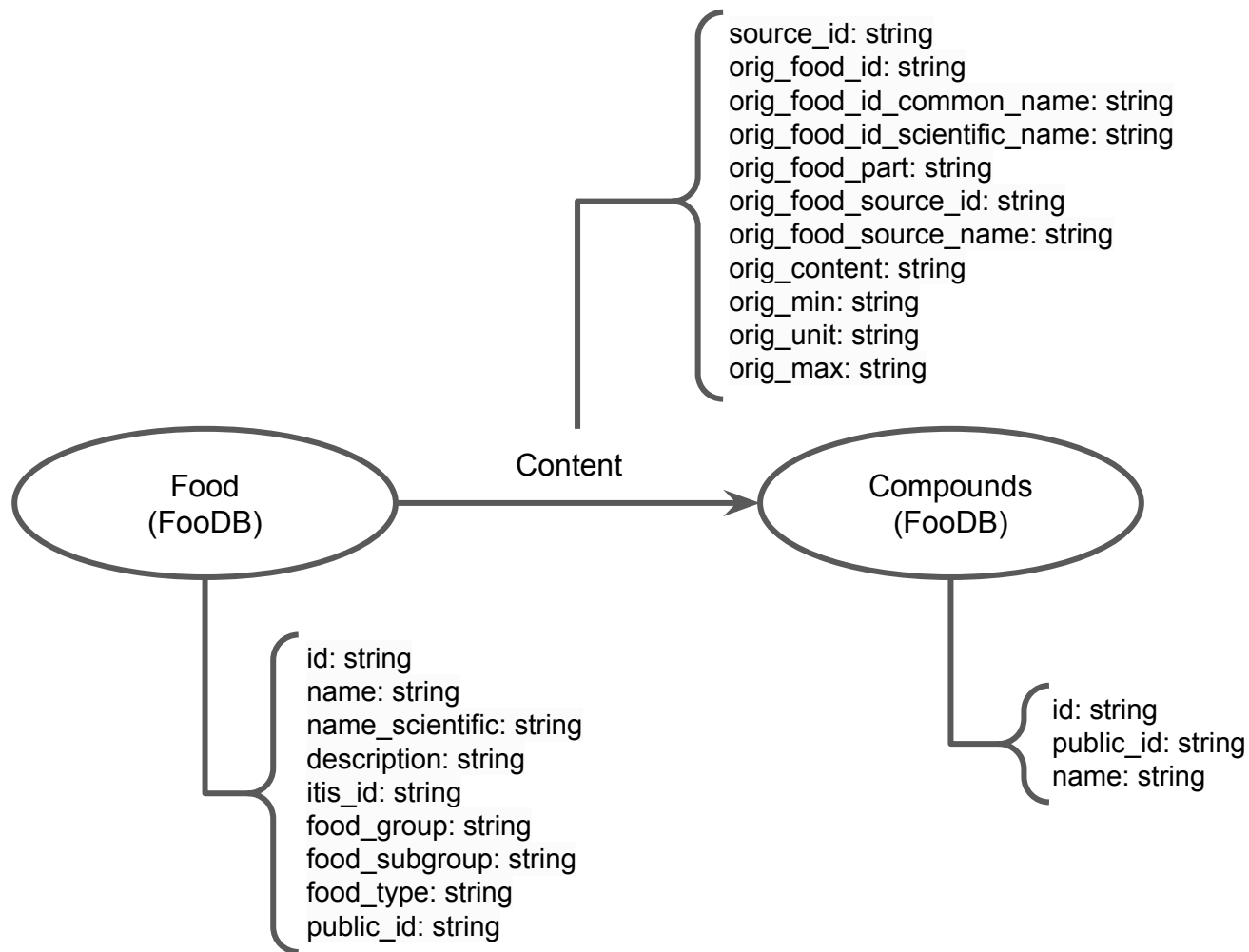
# Usar TF-IDF para converter os nomes em vetores
vectorizer = TfidfVectorizer()
tfidf_matrix = vectorizer.fit_transform(fcid_data['processed'].tolist() + fooddb_data['processed'].tolist())

# Calcular a similaridade de cosseno
cosine_sim = cosine_similarity(tfidf_matrix[:len(fcid_data)], tfidf_matrix[len(fcid_data):])

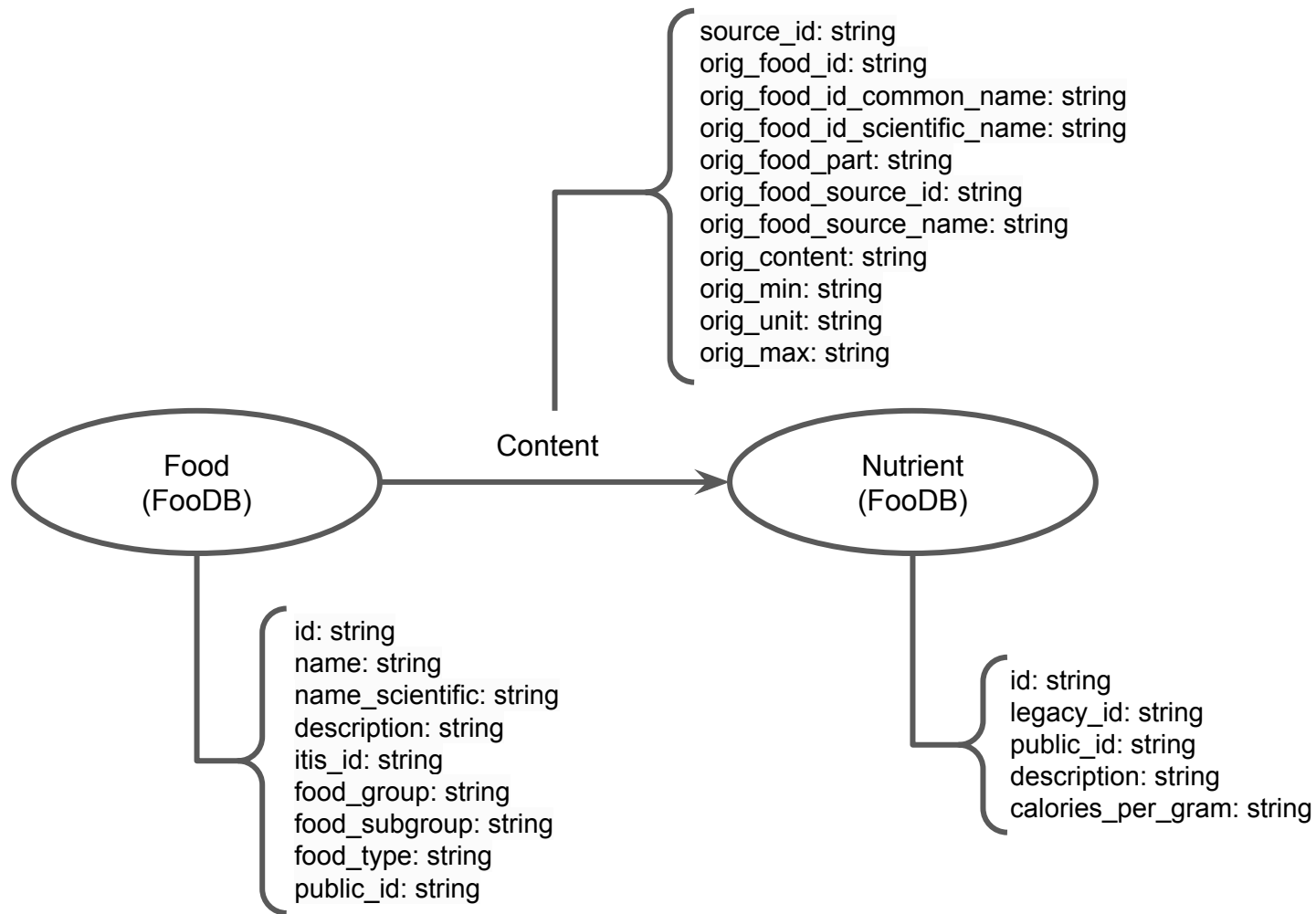
# Encontrar a melhor correspondência para cada alimento do FCID
matches = []
for idx, row in enumerate(cosine_sim):
    best_match_idx = row.argmax()
    fcid_food = fcid_data.iloc[idx]['FCID_Desc']
    fooddb_food = fooddb_data.iloc[best_match_idx]['name']
    matches.append((fcid_food, fooddb_food))

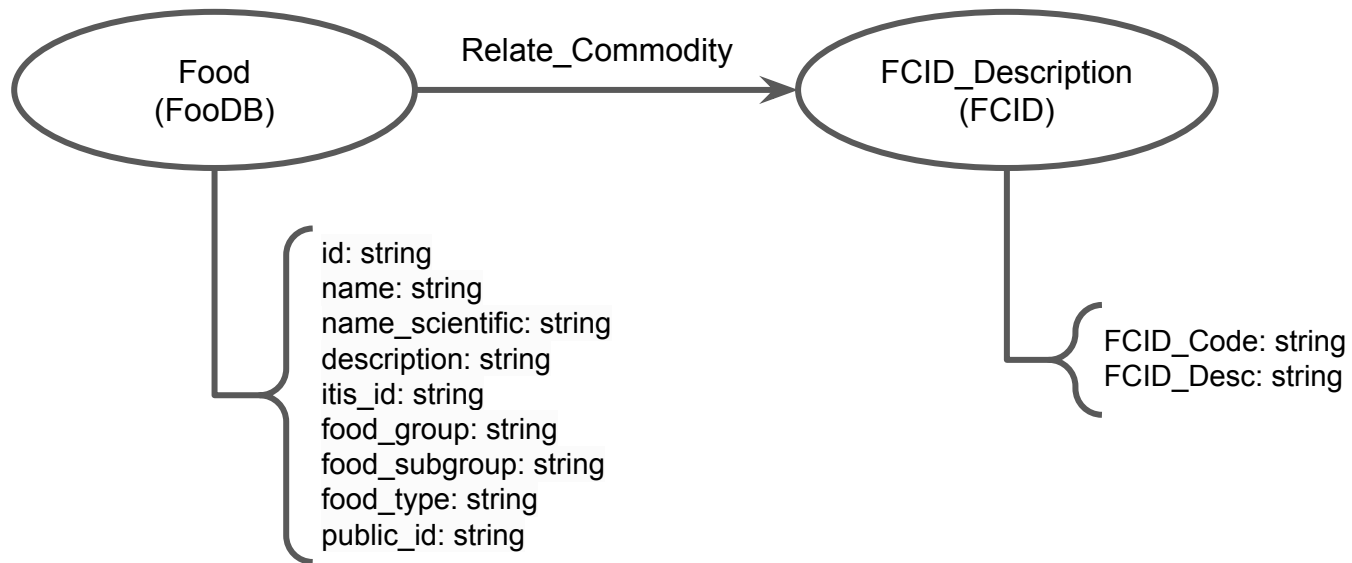
# Criar um DataFrame com as correspondências
matches_df = pd.DataFrame(matches, columns=['FCID_Food', 'FoodDB_Food'])

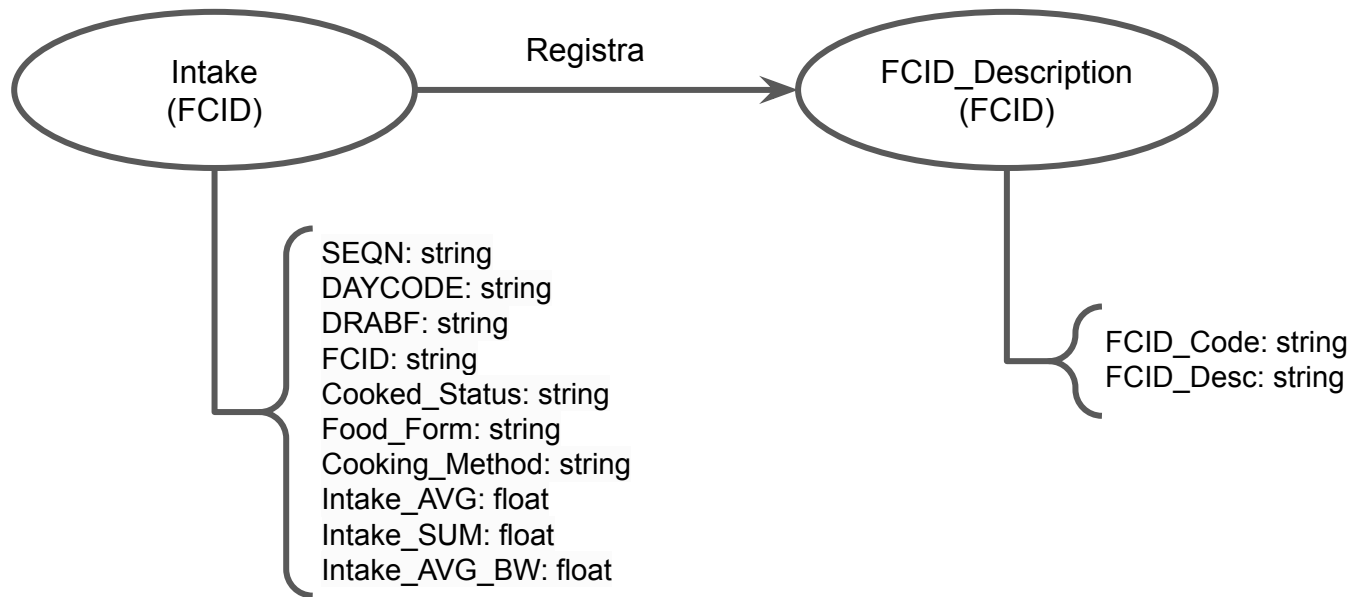
# Salvar o DataFrame em um arquivo CSV
matches_df.to_csv('bases/relação/food_matches.csv', index=False)
```

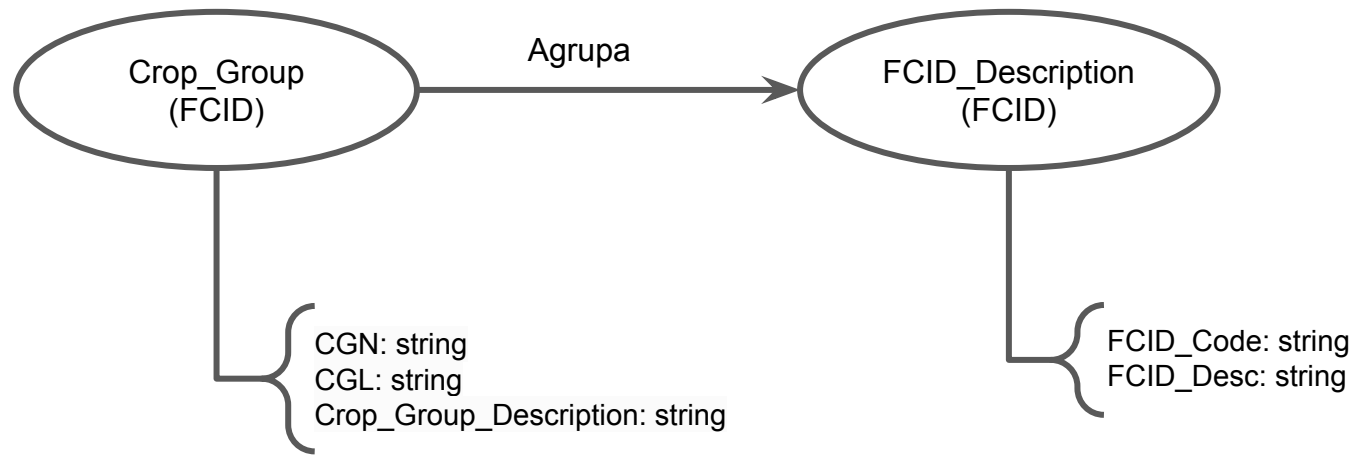


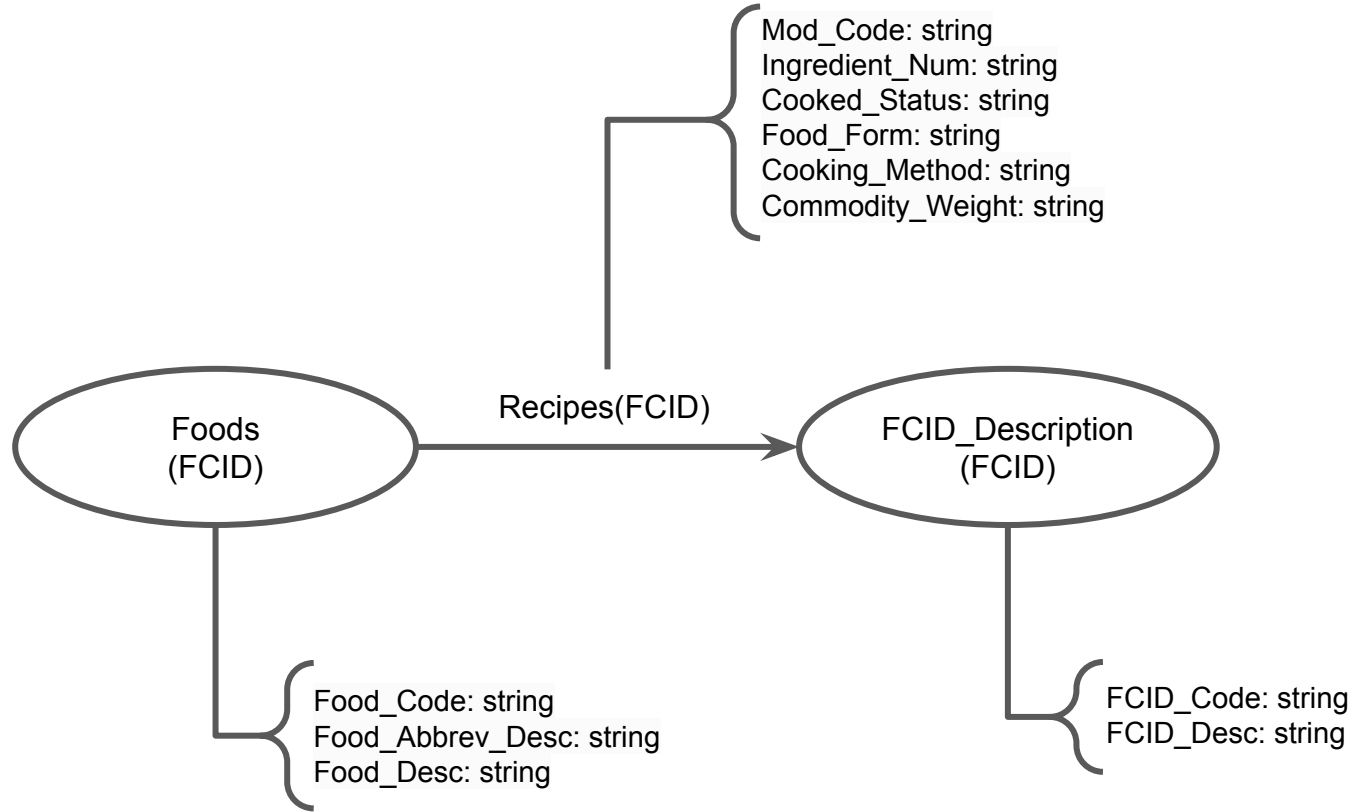


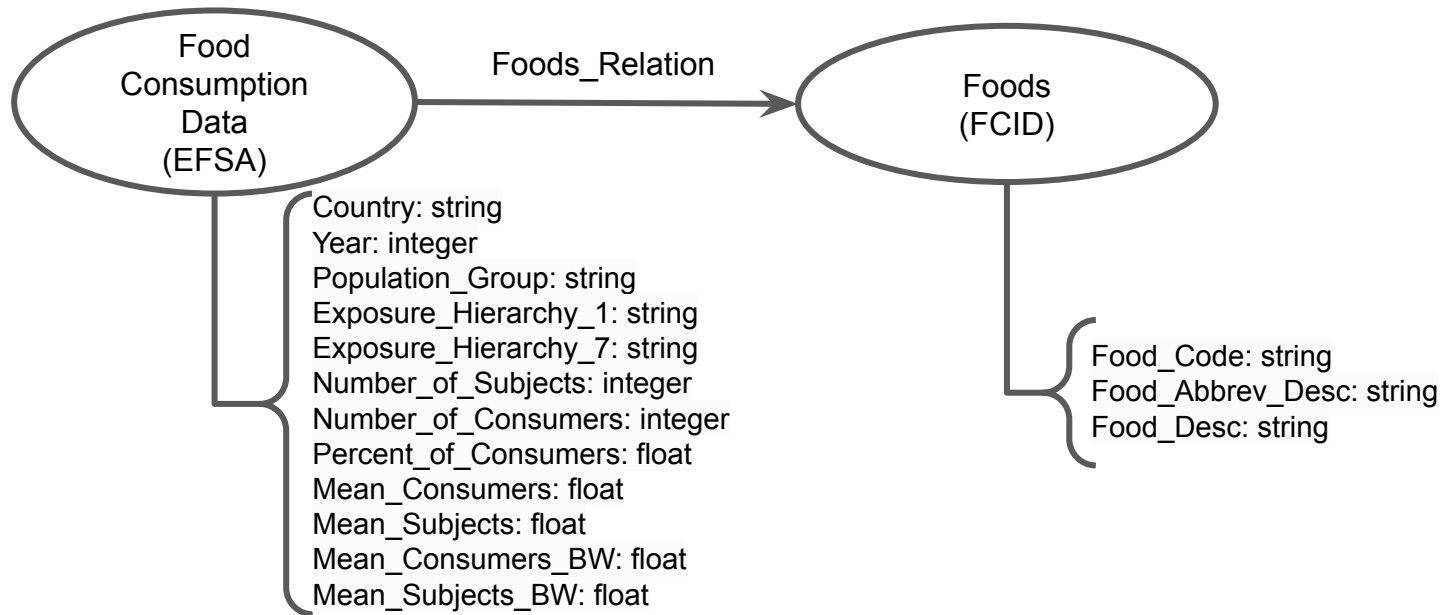








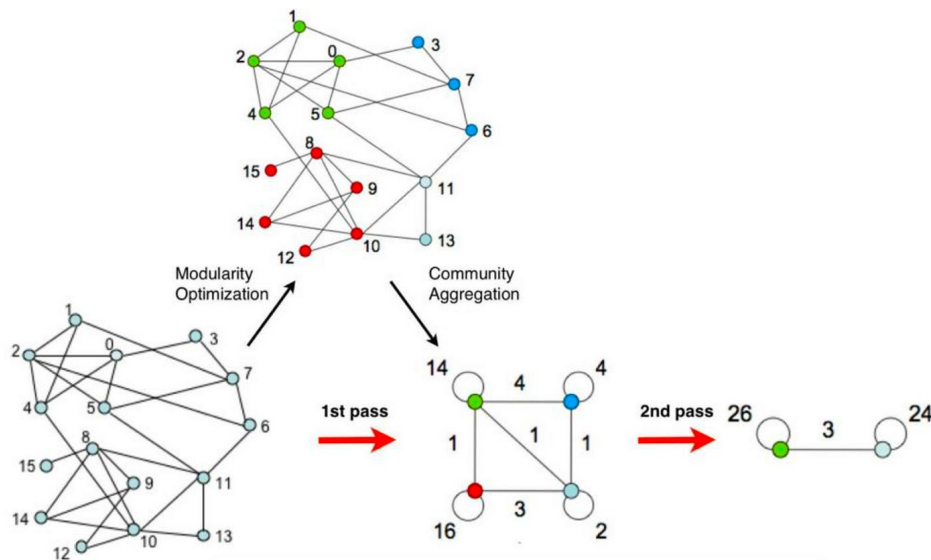




## Perguntas Implementadas (Modelo de Grafos)

Pergunta 1) No contexto das receitas de alimentos, é possível localizarmos uma comunidade de ingredientes que aparecem juntos em várias receitas?

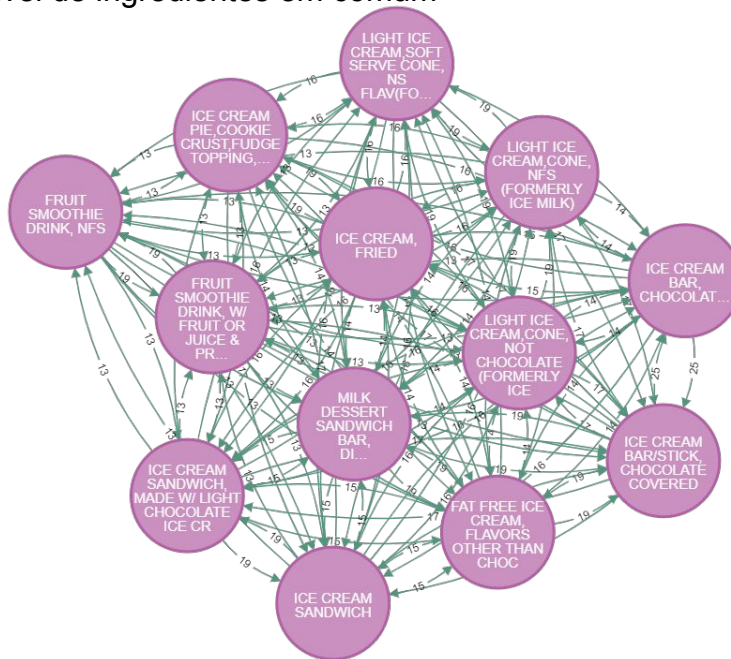
Começamos ligando 2 nós que representam ingredientes sempre que eles estão em uma mesma receita, a fim de aproximar o conceito de 2 ingredientes ligados à ideia de comunidade da pergunta. Utilizamos então, o Algoritmo de Louvain para detectar comunidades nessa nova projeção de grafo:



# Perguntas Implementadas (Modelo de Grafos)

Pergunta 2) Dado que uma pessoa gosta de um alimento, é possível prever que ela gostará de um outro alimento?

Se uma pessoa gosta de um alimento, certamente é por causa dos ingredientes que estão na receita daquele alimento, e assim ela provavelmente gostaria de um alimento suficientemente similar ao primeiro. Vamos ligar 2 alimentos se suas receitas possuem um número razoável de ingredientes em comum

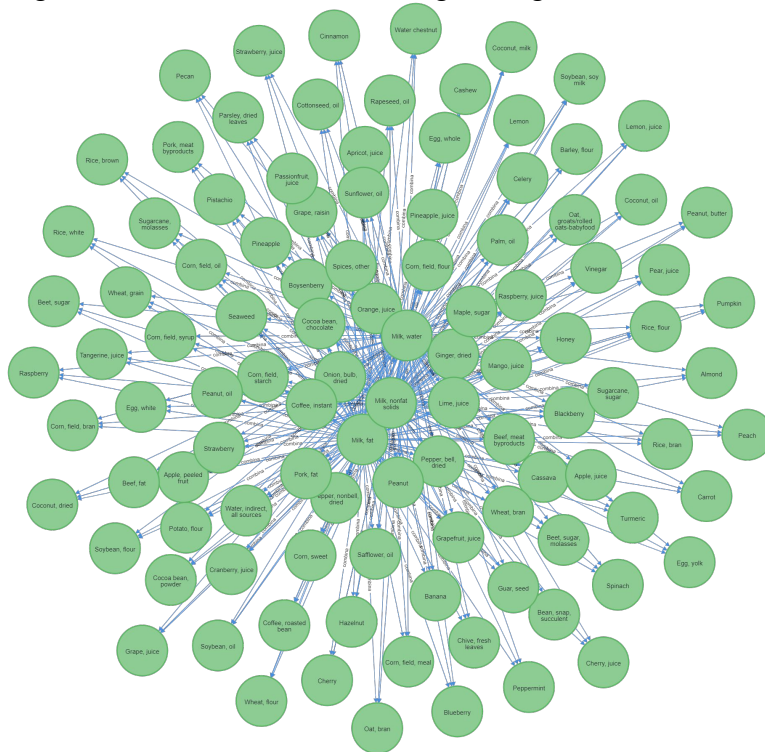




## Perguntas Implementadas (Modelo de Grafos)

Pergunta 3) Qual o ingrediente que mais combina com outros?

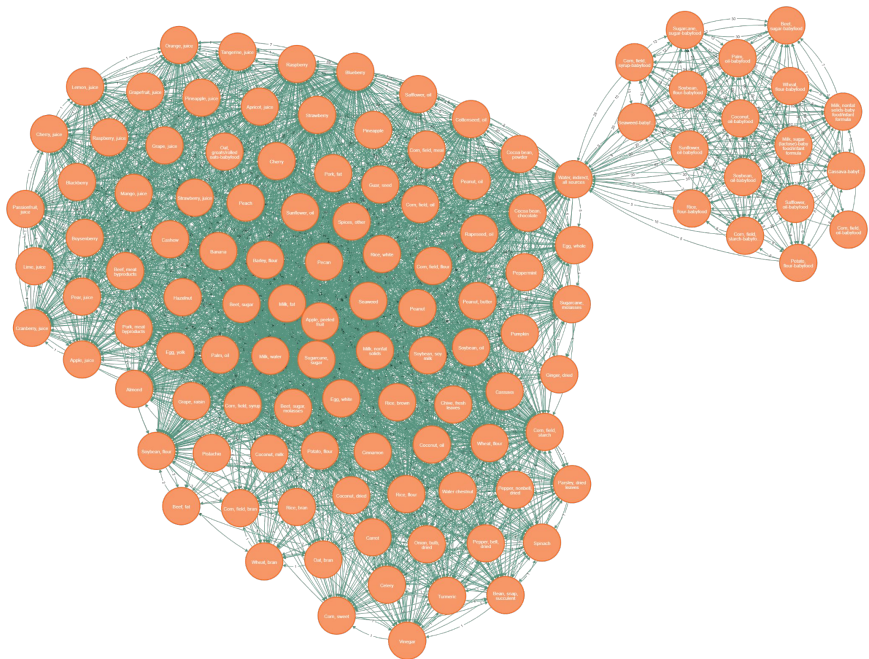
Usamos o conceito de centralidade de grau para detectar ingredientes que estão juntos com outros em muitas receitas. Achamos 3 nós com centralidade de grau = 105, fornecendo o seguinte grafo isolado



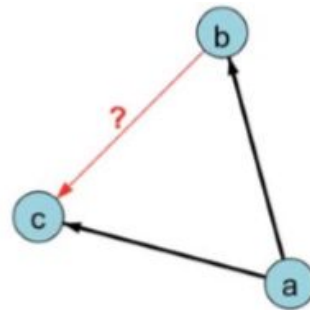
# Perguntas Implementadas (Modelo de Grafos)

Pergunta 4) Como se dá a clusterização dos ingredientes, considerando a relação de combinação entre cada par de ingredientes em uma mesma receita?

Analizamos o coeficiente de clusterização local para cada nó, escolhendo apenas aqueles que têm esse coeficiente igual a 1.0 (o máximo), obtivemos um grafo que acaba retratando 2 grupos de comunidades que quase não se intersectam, como segue:



$$C_i := \frac{2 \cdot k(i)}{d_i(d_i - 1)} = \frac{k(i)}{\binom{d_i}{2}}$$



# Perguntas não Implementadas (Modelo Relacional)

Pergunta 1)

Como o consumo médio na Europa mudou ao longo dos anos de 1997 até 2019?

Pergunta 2)

Quais os alimentos que mais contribuem para o total calórico de cada receita?

Pergunta 3)

Quais os alimentos que cada grupo populacional mais consome? E qual o perfil calórico desses alimentos?

Pergunta 4)

Para cada grupo de alimento, qual o alimento pertencente a este grupo que é mais consumido?

Pergunta 5)

Qual o alimento mais consumido em cada país, levando em conta o peso do consumidor?

Também, como seria o perfil calórico desses alimentos?

Pergunta 6)

Qual grupo de alimentos é mais utilizado para fazer receitas?