# Database Docs

## EER Diagram

# EER to Relational Mapping
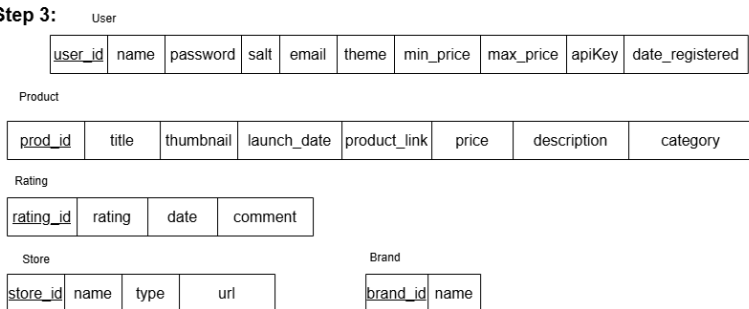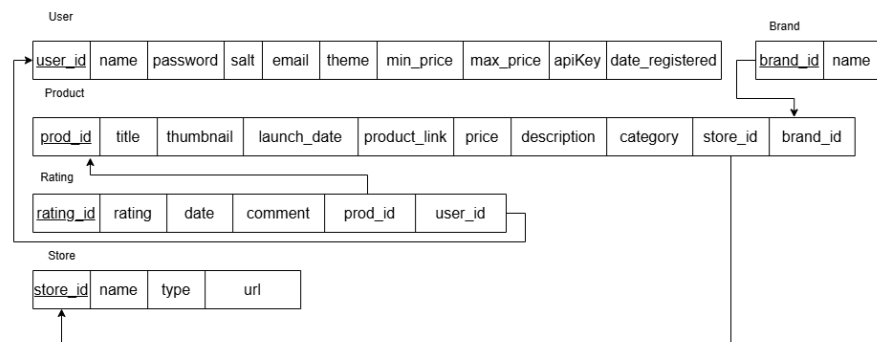
## Step 1:

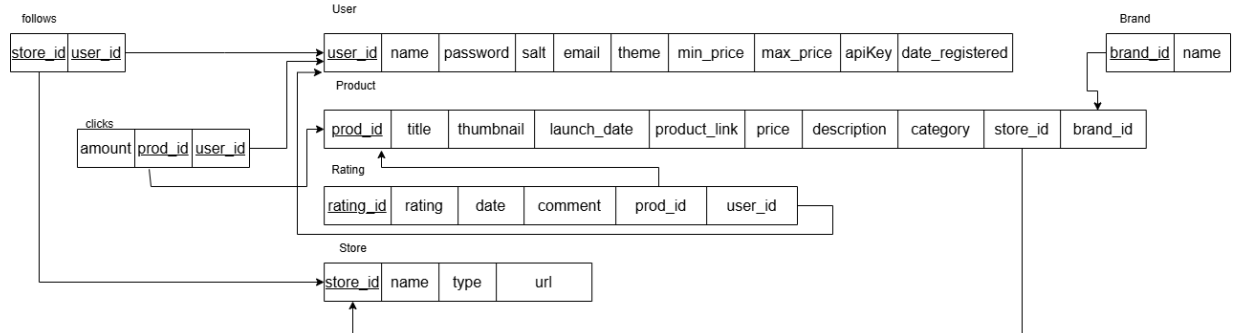**User**

| user_id | name | password | salt | email | theme | min_price | max_price | apiKey | date_registered |
|---------|------|----------|------|-------|-------|-----------|-----------|--------|-----------------|

**Product**

| prod_id | title | thumbnail | launch_date | product_link | price | description | category |
|---------|-------|-----------|-------------|--------------|-------|-------------|----------|

**Rating**

| rating_id | rating | date | comment |
|-----------|--------|------|---------|

**Store**

| store_id | name | type | url |
|----------|------|------|-----|

**Brand**

| brand_id | name |
|----------|------|

## Step 2:

N/A

## Step 3:

**User**

| user_id | name | password | salt | email | theme | min_price | max_price | apiKey | date_registered |
|---------|------|----------|------|-------|-------|-----------|-----------|--------|-----------------|

**Product**

| prod_id | title | thumbnail | launch_date | product_link | price | description | category |
|---------|-------|-----------|-------------|--------------|-------|-------------|----------|

**Rating**

| rating_id | rating | date | comment |
|-----------|--------|------|---------|

**Store**

| store_id | name | type | url |
|----------|------|------|-----|

**Brand**

| brand_id | name |
|----------|------|

## Step 4:

**User**

| user_id | name | password | salt | email | theme | min_price | max_price | apiKey | date_registered |
|---------|------|----------|------|-------|-------|-----------|-----------|--------|-----------------|

**Brand**

| brand_id | name |
|----------|------|

**Product**

| prod_id | title | thumbnail | launch_date | product_link | price | description | category | store_id | brand_id |
|---------|-------|-----------|-------------|--------------|-------|-------------|----------|----------|----------|

**Rating**

| rating_id | rating | date | comment | prod_id | user_id |
|-----------|--------|------|---------|---------|---------|

**Store**

| store_id | name | type | url |
|----------|------|------|-----|

## Step 5:

**follows**

| store_id | user_id |
|----------|---------|

**User**

| user_id | name | password | salt | email | theme | min_price | max_price | apiKey | date_registered |
|---------|------|----------|------|-------|-------|-----------|-----------|--------|-----------------|

**Brand**

| brand_id | name |
|----------|------|

**clicks**

| amount | prod_id | user_id |
|--------|---------|---------|

**Product**

| prod_id | title | thumbnail | launch_date | product_link | price | description | category | store_id | brand_id |
|---------|-------|-----------|-------------|--------------|-------|-------------|----------|----------|----------|

**Rating**

| rating_id | rating | date | comment | prod_id | user_id |
|-----------|--------|------|---------|---------|---------|

**Store**

| store_id | name | type | url |
|----------|------|------|-----|

**Step 6: NOT APPLICABLE**

**Step 7: NOT APPLICABLE**

**Step 8:**

**follows**

| store_id | user_id |
|----------|---------|

**User**

| user_id | name | password | salt | email | theme | min_price | max_price | apiKey | date_registered |
|---------|------|----------|------|-------|-------|-----------|-----------|--------|-----------------|

**store_Owner**

| user_id | registrationNo | store_id |
|---------|----------------|----------|

**customer**

| user_id |
|---------|

**Brand**

| brand_id | name |
|----------|------|

**clicks**

| amount | prod_id | user_id |
|--------|---------|---------|

**Product**

| prod_id | title | thumbnail | launch_date | product_link | price | description | category | store_id | brand_id |
|---------|-------|-----------|-------------|--------------|-------|-------------|----------|----------|----------|

**Rating**

| rating_id | rating | date | comment | prod_id | user_id |
|-----------|--------|------|---------|---------|---------|

**Store**

| store_id | name | type | url |
|----------|------|------|-----|

**Step 9: NOT APPLICABLE**

**FINAL**

**follows**

| store_id | user_id |
|----------|---------|

**User**

| user_id | name | password | salt | email | theme | min_price | max_price | apiKey | date_registered |
|---------|------|----------|------|-------|-------|-----------|-----------|--------|-----------------|

**store_Owner**

| user_id | registrationNo | store_id |
|---------|----------------|----------|

**customer**

| user_id |
|---------|

**Brand**

| brand_id | name |
|----------|------|

**clicks**

| amount | prod_id | user_id |
|--------|---------|---------|

**Product**

| prod_id | title | thumbnail | launch_date | product_link | price | description | category | store_id | brand_id |
|---------|-------|-----------|-------------|--------------|-------|-------------|----------|----------|----------|

**Rating**

| rating_id | rating | date | comment | prod_id | user_id |
|-----------|--------|------|---------|---------|---------|

**Store**

| store_id | name | type | url |
|----------|------|------|-----|

# Relational Schema



**ratings**
- 🔑 rating_id INT
- ◇ rating DECIMAL(2,1)
- ◇ date DATE
- ◇ comment TEXT
- ◇ product_id INT
- ◇ user_id_ratings INT
- ◆ users_user_id INT
- Indexes ▶

**products**
- 🔑 product_id INT
- ◇ title VARCHAR(255)
- ◇ thumbnail VARCHAR(255)
- ◇ launch_date DATE
- ◇ product_link VARCHAR(25...
- ◇ price DECIMAL(10,2)
- ◇ description TEXT
- ◇ category VARCHAR(100)
- ◇ store_id INT
- ◇ brand_id INT
- Indexes ▶

**brand**
- 🔑 brand_id INT
- ◇ name VARCHAR(255)
- Indexes ▶

**users**
- 🔑 user_id INT
- ◇ name VARCHAR(255)
- ◇ password VARCHAR(255)
- ◇ salt VARCHAR(255)
- ◇ email VARCHAR(255)
- ◇ theme ENUM(...)
- ◇ min_price DECIMAL(10,2)
- ◇ max_price DECIMAL(10,...
- ◇ apikey VARCHAR(255)
- ◇ user_type ENUM(...)
- ◇ date_registered DATE
- Indexes ▶

**clicks**
- 🔑 user_id INT
- 🔑 product_id I...
- ◇ amount INT
- Indexes ▶

**store**
- 🔑 store_id INT
- ◇ name VARCHAR(255)
- ◇ type VARCHAR(100)
- ◇ url VARCHAR(255)
- Indexes ▶

**follo...**
- 🔑 user_id INT
- 🔑 store_id INT
- Indexes ▶

**store_owner**
- 🔑 user_id INT
- ◇ registration_no VARCHAR(25...
- ◇ store_id INT
- Indexes ▶

**custom...**
- 🔑 user_id INT
- Indexes ▶

# Data Population Method

We mainly used serpapi to get product data. We kept the data to South Africa to keep the currency consistent for South African stores.

Here is an example of a get request in the serpapi playground using google shopping:

https://serpapi.com/playground?engine=google_shopping&q=stationary &location=South+Africa&google_domain=google.co.za&gl=za&hl=en

This returns a JSON object with a lot of parameters. The usable parameters are under shopping_results.

The parameters we will be retrieving is title, extracted_price, product_link, thumbnail and source.

The data is all stored in a JSON file according to the search term used to get the data and then we used a javascript file to retrieve the data and make an api call to insert the data into the database. For the category attribute we set the category to the search term used to retrieve the data from the api.

Before adding the products though we used the source to generate all the stores needed in the database and used ChatGPT to add links to all the stores and set the store types in the database.

When adding the products we then link the source attribute by using the store_id in the products table.

All other data items were made by using the website as we were testing.

# Query Optimisation

The query we will be optimizing is the following:

SELECT p.*, AVG(r.rating) as average_rating

    FROM products p

    LEFT JOIN ratings r ON p.product_id = r.product_id

    GROUP BY p.product_id ORDER BY average_rating DESC

Currently it takes 0.078 sec to retrieve the data:

| 283 row(s) returned | 0.078 sec / 0.000 sec |
|---|---|

In relational algebra the query can be written as follows:

$\pi$ (p.product_id, title, thumbnail, launch_date, product_link, price, description, category, store_id, brand_id, average_rating) ($\sigma$(p.product_id = r.product_id OR r.product_id = NULL (products x

(product_id $\mathcal{F}$(AVERAGE(rating) AS average_rating)(ratings)))

After optimising the relational algebra, we get:

$\pi$ (p.product_id, title, thumbnail, launch_date, product_link, price, description, category, store_id, brand_id, average_rating) (products $\bowtie$product_id=prod_id $\rho$(average_rating, prod_id) ($\pi$(AVERAGE(rating), product_id)(product_id $\mathcal{F}$(AVERAGE(rating))(ratings))))

Implementing this into SQL yields the following:

SELECT p.*, average_rating

    FROM products AS p

    LEFT JOIN (

SELECT AVG(r.rating) AS average_rating, product_id FROM ratings AS r GROUP BY r.product_id

) AS r ON p.product_id = r.product_id

        ORDER BY average_rating DESC;

This now takes 0.016 sec to retrieve the data

| 283 row(s) returned | 0.016 sec / 0.000 sec |
|---|---|

This means that we have speed up the retrieval process by about 5 times!

This is because we are doing the aggregation first and then joining the tables.