

# **Report - DataBase Project - IFEBY140**

# Table of contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>3</b>  |
| <b>2</b> | <b>Analyzis</b>  | <b>4</b>  |
| 2.1      | Redudant Attributes . . . . .                                  | 4         |
| 2.2      | First Normal Form . . . . .                                    | 4         |
| 2.3      | Functional dependencies : . . . . .                            | 5         |
| 2.3.1    | Non used functional dependencies . . . . .                     | 5         |
| 2.4      | Miscellaneous . . . . .  | 5         |
| 2.5      | Constraints . . . . .  | 5         |
| 2.6      | Triggers . . . . .   | 6         |
| 2.7      | Functional dependencies after decomposition . . . . .          | 6         |
| 2.7.1    | Boyce-Codd verification . . . . .                              | 7         |
| 2.8      | Indexes . . . . .  | 7         |
| <b>3</b> | <b>Modelization</b>  | <b>8</b>  |
| 3.1      | ER schema . . . . .  | 8         |
| 3.2      | Tables . . . . .   | 8         |
| <b>4</b> | <b>Implementation</b>  | <b>9</b>  |
| 4.1      | Execution . . . . .  | 9         |
| 4.2      | Code Details . . . . .   | 9         |
| 4.2.1    | Populate Transports . . . . .                                  | 9         |
| 4.2.2    | Populate Sub-Events . . . . .                                  | 10        |
| <b>5</b> | <b>Request examples</b>  | <b>11</b> |
| 5.1      | Outputs . . . . .  | 11        |
| 5.1.1    | Statistics . . . . .   | 11        |
| 5.1.2    | The 10 station with the most events . . . . .                  | 11        |
| 5.1.3    | The title of the 10 events with the most occurrences . . . . . | 11        |
| 5.1.4    | The 10 most used tag . . . . .                                 | 12        |
| 5.1.5    | The 10 biggest group . . . . .                                 | 12        |

# 1 Introduction

**Que Faire à Paris ?** is a french website that list events happing in Paris.

The event are archived at this [address](#) in multiple formats including csv.

In this project we propose a normalized database to store the events based on the csv archive.

Course constraint :

- Code everything in psql.

- Don't use pgsql functions to extract the data.

## 2 Analyzis

Initial Columns :

```
['ID' 'URL' 'Titre' 'Chapeau' 'Description' 'Date de début' 'Date de fin'
 'Occurrences' 'Description de la date' "URL de l'image"
 "Texte alternatif de l'image" "Crédit de l'image" 'Mots clés'
 'Nom du lieu' 'Adresse du lieu' 'Code postal' 'Ville'
 'Coordonnées géographiques' 'Accès PMR' 'Accès mal voyant'
 'Accès mal entendant' 'Transport' 'Url de contact' 'Téléphone de contact'
 'Email de contact' 'URL Facebook associée' 'URL Twitter associée'
 'Type de prix' 'Détail du prix' "Type d'accès" 'URL de réservation'
 'URL de réservation - Texte' 'Date de mise à jour' 'Image de couverture'
 'Programmes' 'En ligne - address_url' 'En ligne - address_url_text'
 'En ligne - address_text' 'title_event' 'audience' 'childrens' 'group']
```

### 2.1 Redudant Attributes

We will remove the following attributes : `description_de_la_date`

### 2.2 First Normal Form

Non atomic attributs :

- List
  - Occurences : separated by '\_' (underscore)
  - Tags : separated by ','
  - Childrens : separated by ','
  - Transport : separated by '\n'
- Multiple attributes
  - Transport : `transport_type`, `transport_line`, `station`, `distance`
  - Geographic\_Coordinates : `longitude`, `latitude`

## 2.3 Functional dependencies :

- $id \rightarrow url$ 
  - For simplicity sake, we consider that we cannot deduce  $id$  from  $url$ . Perhaps, we will use this property to identify  $sub\_events\_id$  in “childrens”
- $id \rightarrow titre, chapeau, description, date\_de\_debut, date\_de\_fin, occurrences, url\_de\_l\_image, texte\_alternatif\_de\_l\_image, credit\_de\_l\_image, mots\_clefs, nom\_du\_lieu, adresse\_du\_lieu, code\_postale, ville, url\_du\_contact, telephone\_de\_contact, email\_de\_contact, url\_facebook\_associee, url\_twitter\_associe, type\_de\_prix, detail\_du\_prix, type\_d\_acces, url\_de\_reservation, url\_de\_reservation\_texte, date\_de\_mise\_a\_jour, image\_de\_couverture, programmes, en\_ligne\_address\_url, en\_ligne\_address\_url\_text, en\_ligne\_address\_text, title\_event, audience, childrens, groupe, transport\_type, transport\_line, station, distance$
- $date\_de\_debut, date\_de\_fin \rightarrow description\_de\_la\_date$ 
  - Removed attribute
- $adresse\_du\_lieu, ville, code\_postale \rightarrow coordonnees\_geographiques$
- $nom\_du\_lieu, adresse\_du\_lieu, ville, code\_postale \rightarrow acces\_pmr, acces\_mal\_voyant, acces\_mal\_entendant$

### 2.3.1 Non used functional dependencies

- In 1NF :  $transport\_station \rightarrow ville$ 
  - It’s true, but for simplicity sake, we’ll ignore it.
- $ville \rightarrow cp$  and  $cp \rightarrow ville$ 
  - It’s neither true in France and Ile de France
- $(titre, date\_de\_debut, date\_de\_fin, nom\_du\_lieu, adresse\_du\_lieu, code\_postale, ville) \rightarrow id$ 
  - We cannot use those attributes as a primary key, because they are often NULL, thus we’ll ignore this functional dependency.

## 2.4 Miscellaneous

We decide to translate everything into English.

## 2.5 Constraints

- Realistic implementation expectations
  - $date\_end > date\_start$
  - NOT NULL :
    - \*  $id, url$
    - \*  $(date\_start \text{ AND } date\_end) \text{ OR } (occurrences)$

- \* title
- Unrealistic implementation expectations
  - url must finish with id
  - all urls must be valid
  - address must be valid
  - phone\_number must be valid
  - every event should have an address, a geographic coordinate and a contact
  - NOT NULL
    - \* a non null lead text, description, date\_start, date\_end, title\_event, price\_type

## 2.6 Triggers

- event\_table.parent\_event\_id : ON DELETE CASCADE
- transport.event\_id : ON DELETE CASCADE

## 2.7 Functional dependencies after decomposition

- Relation : **geographic\_correspondance**
  - address\_street, address\_zipcode, address\_city → lat, lon
- Relation : **address\_table**
  - address\_name, address\_street, address\_zipcode, address\_city → pmr, blind, deaf
- Relation : **event\_table**
  - event\_id → event\_url, title, lead\_text, event\_description, date\_start, date\_end, cover\_url, cover\_alt, cover\_credit, address\_name, address\_street, address\_zipcode, address\_city, price\_type, price\_detail, access\_type, access\_link, access\_link\_text, updated\_at, image\_couverture, programs, title\_event, audience, contact\_url, contact\_phone, contact\_mail, contact\_facebook, contact\_twitter, address\_url, address\_url\_text, address\_text, keyword, group\_name, parent\_event\_id
- Relation : **occurence**
  - event\_id, occurence\_date (no functional dependancy)
- Relation : **tag**
  - event\_id, keyword (no functional dependancy)
- Relation : **transport**
  - event\_id, transport\_type, transport\_line, station → distance

### 2.7.1 Boyce-Codd verification

- 1NF : Atomicity and no list attributes
  - Checked
- 2NF : Attributes cannot depend only on a subpart of the primary key
  - Checked
- 3NF : No transitive dependances
  - Checked
- Boyce-Codd :  $X \rightarrow Y \implies X$  is a superkey keyword
  - Checked

## 2.8 Indexes

We often want to identify events using :

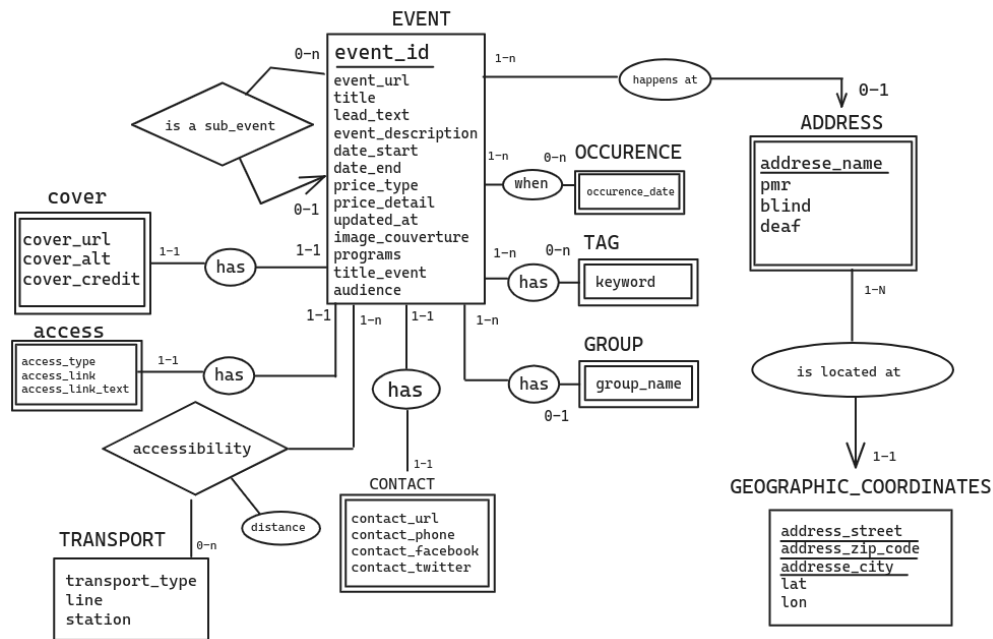
- the title
  - event\_table (title)
  - event\_table (title\_event)
- the period
  - event\_table (date\_begin)
  - event\_table (date\_end)
- the arrondissement of Paris
  - event\_table (address\_zipcode)
- using the url
  - event\_table (event\_url)
- event\_table (parent\_event\_id)

We also want to do agregations on keywords :

- tag (keyword)

## 3 Modelization

### 3.1 ER schema



The following decomposition, satisfies Boyce-codd normal form.

### 3.2 Tables

**geographic\_correspondance**(address\_street, address\_zipcode, address\_city | lat, lon)

**address\_table**(address\_name, address\_street, address\_zipcode, address\_city | pmr, blind, deaf)

**event\_table** (event\_id | event\_url, title, lead\_text, event\_description, date\_start, date\_end, cover\_url, cover\_alt, cover\_credit, address\_name, address\_street, address\_zipcode, address\_city, price\_type, price\_detail, access\_type, access\_link, access\_link\_text, updated\_at, image\_couverture, programs, title\_event, audience, contact\_url, contact\_phone, contact\_mail, contact\_facebook, contact\_twitter, address\_url, address\_url\_text, address\_text, keyword, group\_name, parent\_event\_id)

**occurence**(event\_id, occurence\_date)

**tag**(event\_id, keyword)

**transport**(event\_id, transport\_type, transport\_line, station | distance)



## 4 Implementation

### 4.1 Execution

Start psql in src folder.

```
cd src
psql -d <database> username
```

Create tables and populate them with :

```
\i XXX_YYY_tables.sql
```

### 4.2 Code Details

We ignore most redundancy errors using `ON CONFLICT DO NOTHING`.

#### 4.2.1 Populate Transports

Raw transport line :

"Métro -> 1 : Bastille (Paris) (304m)

Bus -> 29618791 : Lyon / Daumesnil - Ledru Rollin (Paris) (246m)

Vélib -> Lacuée - Lyon (121.91m)

`<a href=""https://www.geovelo.fr/paris/route?to=2.370458543300628,48.849268481958404"">Calculer`

1. List → Row

- Split over line jump using : `unnest(string_to_array(transport, E'\n')) AS transport`

2. Remove unnecessary HTML content

- `WHERE transport NOT LIKE '<%'`

3. Extract (transport\_type, transport\_line, station + distance)

- use `split_part` over -> and :

4. Vélib case

- If station is empty, it means it is a Vélib, and the station was wrongly splitted into transport\_line. Then we swap `station` and `transport_line`

## 5. Extract distance

- It can be done using a combination of `substring`, `reverse` and `length` functions.
- Distance information is always between the last occurrence of ( and m).

### 4.2.2 Populate Sub-Events

Raw sub event line example:

"Rencontre avec Jørn Lier Horst (<https://www.paris.fr/evenements/rencontre-avec-jorn-lier-horst>)

#### 1. List → Row

- Split over ; delimiters using : `unnest(string_to_array(children, ';'))`

#### 2. Extract id

- The id is always at the end of the url. It can be extracted knowing it's always between the last - and ) of the string.

## 5 Request examples

See the interactive preview of the data with :

```
\i XXX_YYY_data.sql
```

### 5.1 Outputs

#### 5.1.1 Statistics

Number of events : 2816

Number of events with pmr, blind or deaf access : 432

Number of events accessibles with “Vélib” : 264

Number of events with multiple transports : 266

#### 5.1.2 The 10 station with the most events

| count | transport_type | transport_line | station                                 |
|-------|----------------|----------------|---|
| 75    | Bus            | 29618791       | Lyon / Daumesnil - Ledru Rollin (Paris) |
| 75    | Métro          | 1              | Bastille (Paris)                        |
| 75    | Vélib          |                | Lacuée - Lyon                           |
| 55    | Vélib          |                | Boyer - Ménilmontant                    |
| 30    | Bus            | 2696           | Pyrénées - Ménilmontant (Paris)         |
| 30    | Métro          | 11             | Jourdain (Paris)                        |
| 29    | Bus            | 397089         | Volontaires - Lecourbe (Paris)          |
| 29    | Vélib          |                | Volontaire - Lecourbe                   |
| 29    | Métro          | 12             | Volontaires (Paris)                     |
| 25    | Métro          | 3              | Gambetta (Paris)                        |

#### 5.1.3 The title of the 10 events with the most occurrences

| count | id    | title  |
|-------|-------|--|
| 9970  | 41529 | Un parcours sonore et urbain   |
| 4031  | 37301 | La Cité de l'Histoire, l'expérience immersive à travers l'Histoire de France |
| 641   | 31363 | Une plongée immersive avec les baleines à l'Aquarium tropical                |

| count | id    | title  |
|-------|-------|--|
| 634   | 11871 | Revivre, les animaux disparus en réalité augmentée                   |
| 572   | 42634 | « Notre-Dame de Paris : au cœur du chantier »                        |
| 552   | 53785 | Montmartre Enchanté insolite : la visite chantée et commentée        |
| 511   | 32056 | Mémorial de l'ancienne gare de déportation de Bobigny                |
| 379   | 32135 | Exposition " Notre-Dame de Paris. Des bâtisseurs aux restaurateurs." |
| 352   | 50933 | Sauvez le cinéma, l'escape game du Grand Rex                         |
| 342   | 31054 | « Félins » au Muséum national d'Histoire naturelle                   |

#### 5.1.4 The 10 most used tag

| count | keyword     |
|-------|-------------|
| 641   | Musique     |
| 623   | Concert     |
| 503   | Atelier     |
| 454   | Enfants     |
| 393   | Loisirs     |
| 312   | Sport       |
| 265   | Expo        |
| 240   | Conférence  |
| 234   | Théâtre     |
| 179   | Littérature |

#### 5.1.5 The 10 biggest group

| count | group_name           |
|-------|----------------------|
| 1184  | Aucun                |
| 634   | Bibliothèques        |
| 602   | Agenda               |
| 214   | Activités DJS        |
| 74    | Centres d'animations |
| 49    | Parcs et jardins     |
| 23    | Nuit Blanche         |
| 22    | Associations         |
| 14    | Musées               |