

## Part B:

### Step 1

I first use the bottom-up method(DO Normalization first) to analyze several provided CSV files to have a rough outline of the overall required structure and the relationships between entities.

*Locations.csv( iso\_code, location\_name, vaccines\_manufacturer(multi-valued), last\_observation\_date(LOD), source\_name, source\_link)*

FDs:

iso\_code → location\_name, LOD, source\_name , source\_link

vaccines\_manufacturer → vaccines\_manufacturer (Manufacturers' Name)

Candidate key:

{ **iso\_code, vaccines\_manufacturer**}<sup>+</sup>={iso\_code, location\_name, LOD, source\_name, source\_website, vaccines\_manufacturer }

Decomposition:

Merges FDs with the same LHS, then form a relation for each FDs

**Locations**(iso\_code, location\_name, LOD, source\_link , source\_name)

**Manufacturers**( Mname )

If no relation contains an entire candidate key of the original relation, add one relation containing key

**Vaccine\_Providers**(iso\_code\*, Mname\*)

*Vaccinations( location, iso\_code, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations\_raw, daily\_vaccinations, total\_vaccinations\_per\_hundred, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, total\_boosters\_per\_hundred, daily\_vaccinations\_per\_million, daily\_people\_vaccinated, daily\_people\_vaccinated\_per\_hundred)*

Since the original relation schema contains too many attributes that are just for statistical purposes, the following are reasons and steps I simplify some unnecessary attributes:

1. I removed per\_hundred attributes because they can be easily calculated by locations' populations and are unnecessary for this database design.
2. I removed daily\_vaccinations\_raw since the GitHub source recommended that any analysis on daily vaccination rates should be conducted using daily\_vaccinations instead.
3. Some locations (countries) didn't report vaccinations regularly, so I will use daily\_vaccinations to calculate the monthly total\_vaccinations.
4. Already know Iso\_code can define location\_name, so I removed it directly here to reduce redundants.

After simplify:

**Records( iso\_code, date, total\_vaccinations, daily\_vaccinations, daily\_people\_vaccinated, people\_fully\_vaccinated, total\_boosters)**

Candidate key:

{ **iso\_code, date** }<sup>+</sup>= {iso\_code, date, total\_vaccinations, daily\_vaccinations, daily\_people\_vaccinated, people\_fully\_vaccinated, total\_boosters }

## Step 2 – Business Rule:

After going through all the csv and analyzing them, I attempted to organize the business rules for this project (based on the GitHub README and my own data analysis):

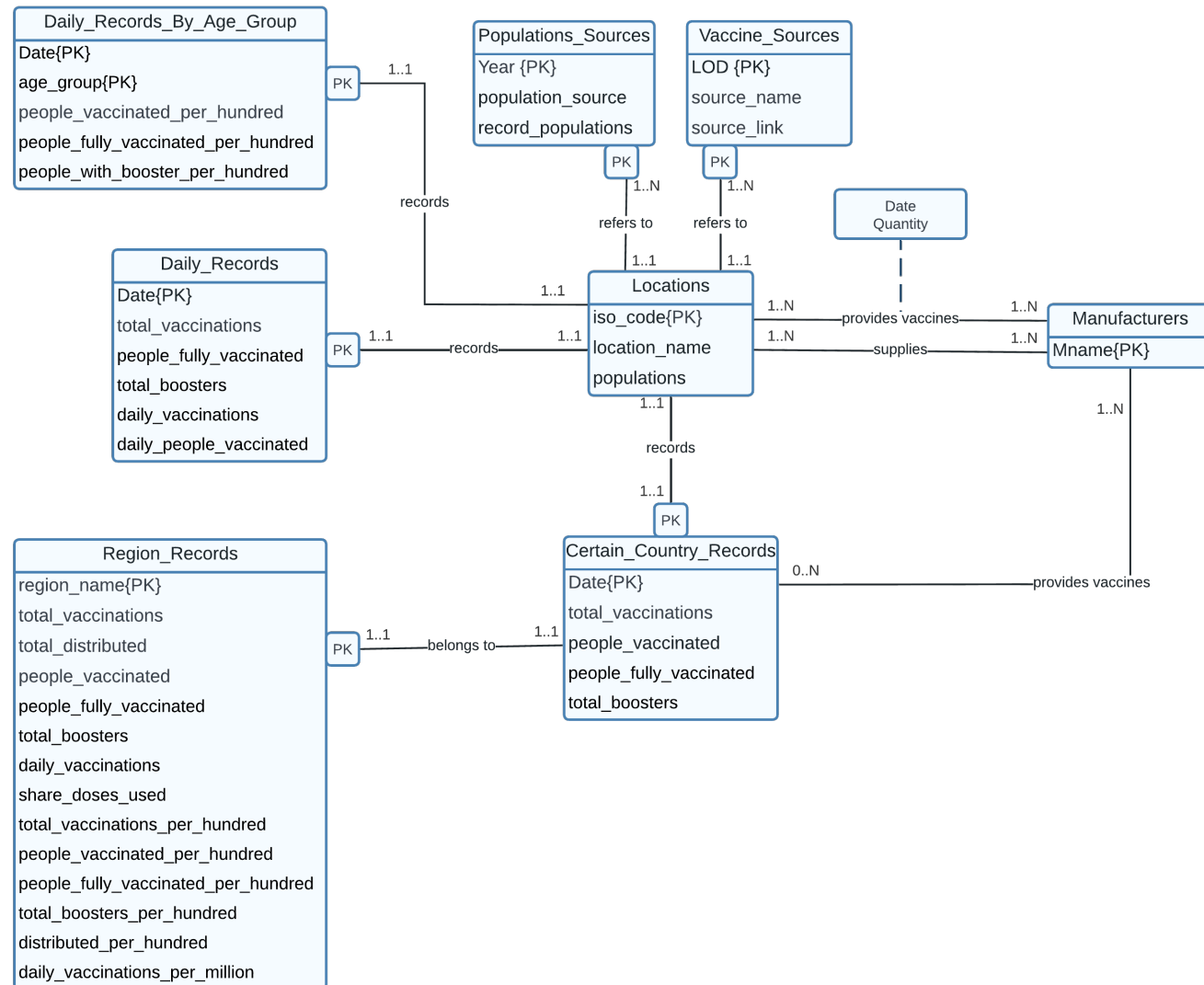
1. **iso\_code** can uniquely identify locations (countries, regions, or continents).
2. **last\_observation\_date (LOD)** is the date of the latest observed data for locations from certain sources.
3. The same **source\_name** may have different **source\_links**, and the same **source\_link** may have different **source\_names**.
4. **Year** is the year of the population data obtained from the **population\_source** for locations.
5. One manufacturer can provide vaccines to multiple locations, and one location can accept vaccines from multiple manufacturers.
6. Manufacturers may not consistently supply vaccines to certain locations.
7. For some regions, manufacturer supply records are only available for broad areas (e.g., European Union or South Africa), while specific country records will have detailed information on which manufacturer provided vaccines to that location on each day.
8. Here, the assumption is that each location records vaccine administration data daily (see the note below), but it should be noted that the recording time span for each location varies (e.g., some may range from 2020/01 to 2023/12, while others may have ended in 2022/01).
9. Although **vaccinations.csv** includes records for almost all locations, certain countries/locations still have their own records in separate CSV files, which contain different types of data compared to **vaccinations.csv**.
10. Some countries also have detailed data for regions/states within the country (e.g., **us\_state\_vaccinations.csv**). Since there is currently no detailed population data for these regions, the **per\_hundred** statistics in this file should be retained.

### Note:

1. **Population** data is additional data introduced (not part of the given dataset for the assignment), which can reduce the need to store various **per\_hundred** statistics from the original data. If needed, calculations can be performed using the population data. The data source is **covid-19-data/scripts/input/un/population\_latest.csv** from the same GitHub repository.
2. Not every region has accurate daily records of vaccine administration data, but due to the use of the 7-day rolling average smoothing method in the original data, dates with missing data are estimated using data from the previous seven days to derive theoretical data.
3. Some locations (countries) reported **people\_fully\_vaccinated** data irregularly, so data calculated on a monthly basis might have differences in sampling dates.

4. Due to the 7-day rolling average smoothing method, using **daily\_vaccinations** to calculate the monthly **total\_vaccinations** may introduce some discrepancies compared to the actual numbers. The same logic applies to **people\_vaccinated** and **daily\_people\_vaccinated**.

### Step 3 – ER Diagram:



## Some assumptions or notes for ER:

1. The one-to-many relationship between locations and sources is to store the different sources from which locations' data is sampled at different times (e.g., the US population data might be sourced from A in 2022 and from B in 2023). If not indicated, the data is taken from the newest sources by default.
2. The relationship between locations and manufacturers for providing vaccines is recorded on a daily basis. The supplies relationship simply records any instance of supply, marking any manufacturer that has ever supplied vaccines to the location as a supplier.
3. Different types of records (such as `by_age_group` or `certain_country_record`) might have different recording periods for the same location, so some corresponding relationships might have a lower limit (participation) of 0.
4. Each location will have only one record per day; there will be no instances of multiple records for the same location on the same day.
5. The **Daily\_Records** still retains the original `total_vaccinations` column to provide accurate actual data when the cumulative total numbers at a certain point in time are needed, rather than summing up the potentially error-prone **daily\_vaccinations** data.
6. The records of the same iso in **Daily\_Records** and **Daily\_Records\_By\_Age\_Group** may have different data collecting timeframes, such as ARG has records data in the whole 2021 in **Daily\_Records\_By\_Age\_Group**, but **Daily\_Records** only has 2021 January records data. Therefore, I treated them separately. We cannot get **Daily\_Records\_By\_Age\_Group** directly from categorizing **Daily\_Records** by age group. The date and the number of `iso_codes` are different in these two csv.

## Step 4 – ER Mapping:

### Step 1: Map Strong Entities

**Locations**(iso\_code, location\_name, populations)

**Manufacturers**( Mname )

### Step 2: Map Weak Entities

**Locations**(iso\_code, location\_name, populations)

**Manufacturers**( Mname )

**Vaccine\_Sources**( iso\_code\*, LOD, source\_name, source\_link)

**Population\_Sources**( iso\_code\*, Year, p\_source\_link, record\_populations)

**Daily\_Records**( iso\_code\*, date, total\_vaccinations, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations, daily\_people\_vaccinated)

**Daily\_Records\_By\_Age\_Group**( iso\_code\*, date, age\_group, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred)

**Certain\_Country\_Records**( iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)

**Region\_Records**( iso\_code\*, date\*, region\_name, total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daly\_vaccinations, share\_doses\_used, total\_vaccinations\_per\_hundred, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, total\_boosters\_per\_hundred, distributed\_per\_hundred, daily\_vaccinations\_per\_million)

Step 3: Map 1:1 Relationships

None

Step 4: Map 1:N Relationships

None

Step 5: Map N:N Relationships

**Provides\_Vaccines**( iso\_code\*, Mname\*, Date, Quantity)

**Suppliers**( iso\_code\*, Mname\*)

**Provides\_Vaccines\_Certain\_Country**(iso\_code\*, Mname\*, Date\*, Sum\_Quantity)

Step 6: Multi-valued Attributes

None

Step 7: Map higher-degree relationships

None

Final Schema:

**Locations**(iso\_code, location\_name, populations)

**Manufacturers**( Mname )

**Vaccine\_Sources**(iso\_code\*, LOD, source\_name, source\_link)

**Population\_Sources**(iso\_code\*, Year, p\_source\_link, record\_populations)

**Daily\_Records**(iso\_code\*, date, total\_vaccinations, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations, daily\_people\_vaccinated)

**Daily\_Records\_By\_Age\_Group**(iso\_code\*, date, age\_group, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred)

**Certain\_Country\_Records**( iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)

**Region\_Records**( iso\_code\*, date\*, region\_name, total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daly\_vaccinations, share\_doses\_used, total\_vaccinations\_per\_hundred, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, total\_boosters\_per\_hundred, distributed\_per\_hundred, daily\_vaccinations\_per\_million)

**Provides\_Vaccines**(iso\_code\*, Mname\*, Date, Quantity)

**Suppliers**(iso\_code\*, Mname\*)

**Provides\_Vaccines\_Certain\_Country**(iso\_code\*, Mname\*, Date\*)

## Step 5- Check Normal Form:

Locations FDs:

iso\_code → location\_name, populations

location\_name → iso\_code, populations

⌘ Using **iso\_code** as the primary key because it is standardized. Although **location\_name** is also unique, it may have spelling inconsistencies.

Manufacturers FDs:

Mname → Mname

Vaccine\_Sources FDs:

iso\_code, LOD → source\_name, source\_link

Population\_Sources FDs:

iso\_code, Year → p\_source\_link, record\_populations

Daily\_Records FDs:



iso\_code, date → total\_vaccinations, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations, daily\_people\_vaccinated

Daily\_Records\_By\_Age\_Group FDs:

iso\_code, date, age\_group → people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred

Certain\_Country\_Records FDs:

iso\_code, date → total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters

Region\_Records FDs:

iso\_code, date, region\_name → total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations, share\_doses\_used, total\_vaccinations\_per\_hundred, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, total\_boosters\_per\_hundred, distributed\_per\_hundred, daily\_vaccinations\_per\_million

Provides\_Vaccines FDs:

iso\_code, Mname, Date → Quantity

Suppliers FDs:

iso\_code → iso\_code

Mname → Mname

Provides\_Vaccines\_Certain\_Country FDs:

iso\_code, Mname, Date → iso\_code, Mname, Date

**Closures:**

Locations: { **iso\_code** }+= { iso\_code, location\_name, populations }

Manufacturers: { **Mname** }+= { Mname }

Vaccine\_Sources: { **iso\_code, LOD** }+= { iso\_code, LOD, source\_name, source\_link }

Population\_Sources: { **iso\_code, Year** }+= { iso\_code, Year, p\_source\_link, record\_populations }

Daily\_Records: { **iso\_code, date** }+= { iso\_code, date, total\_vaccinations, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations, daily\_people\_vaccinated }

Daily\_Records\_By\_Age\_Group: { **iso\_code, date, age\_group** }+= { iso\_code, date, age\_group, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred }

Certain\_Country\_Records: { **iso\_code, date** }+= { iso\_code, date, t total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters }

Region\_Records: { **iso\_code, date, region\_name** }+= { iso\_code, date, region\_name, total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daly\_vaccinations, share\_doses\_used, total\_vaccinations\_per\_hundred, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, total\_boosters\_per\_hundred, distributed\_per\_hundred, daily\_vaccinations\_per\_million }

Provides\_Vaccines: { **iso\_code, Mname, Date** }+= { iso\_code, Mname, Date, Quantity }

Suppliers: { **iso\_code, Mname** }+= { iso\_code, Mname }

Provides\_Vaccines\_Certain\_Country: { **iso\_code, Mname, Date** }+= { iso\_code, Mname, Date }

## Highest NF:

No Duplicated row/data or multivalued → pass 1NF

No partial dependency on non-prime attributes → pass 2NF

No transitive dependency on non-prime attributes → pass 3NF

Every schema is at 3NF.

### Final Schemas:

**Locations**(iso\_code, location\_name, populations)

**Manufacturers**( Mname )

**Vaccine\_Sources**( iso\_code\*, LOD, source\_name, source\_link)

**Population\_Sources**( iso\_code\*, Year, p\_source\_link, record\_populations)

**Daily\_Records**( iso\_code\*, date, total\_vaccinations, people\_fully\_vaccinated, total\_boosters, daily\_vaccinations, daily\_people\_vaccinated)

**Daily\_Records\_By\_Age\_Group**( iso\_code\*, date, age\_group, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, people\_with\_booster\_per\_hundred)

**Certain\_Country\_Records**( iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)

**Region\_Records**( iso\_code\*, date\*, region\_name, total\_vaccinations, total\_distributed, people\_vaccinated, people\_fully\_vaccinated, total\_boosters, daly\_vaccinations, share\_doses\_used, total\_vaccinations\_per\_hundred, people\_vaccinated\_per\_hundred, people\_fully\_vaccinated\_per\_hundred, total\_boosters\_per\_hundred, distributed\_per\_hundred, daily\_vaccinations\_per\_million)

**Provides\_Vaccines**( iso\_code\*, Mname\*, Date, Quantity)

**Suppliers**( iso\_code\*, Mname\*)

**Provides\_Vaccines\_Certain\_Country**(iso\_code\*, Mname\*, Date\*)

※ **Note for modifying the CSV data:**

1. The original **vaccinations\_by\_manufacturer.csv** (later called **Provides\_vaccines.csv**) only has **location\_name**. Therefore, I used Excel functions to convert them to the corresponding **iso\_codes** in **locations.csv**. Also, in the last section of the table, the **location\_name** was listed as European Union, which was not included in the original **locations.csv** and did not have a corresponding **iso\_code**.
2. In **vaccinations.csv**, there are additional **iso\_codes**: OWID\_AFR, OWID\_ASI, OWID\_EUN, OWID\_EUR, OWID\_HIC, OWID\_LIC, OWID\_LMC, OWID\_NAM, OWID\_OCE, OWID\_SAM, OWID\_UMC, and OWID\_WRL, including the **iso\_code** for the European Union. Therefore, I added them to **locations.csv** and also replaced European Union in **Provides\_vaccines.csv** with the corresponding **iso\_code** to correctly import the data into the database without triggering the Foreign Key Constraint.
3. If specific data is 0 or NULL, it means that the data was not obtained and does not necessarily mean that the value is 0. For example, in the **daily\_record**, some locations have data marked as 0 for 2020, which does not mean that there were no vaccinations in that country in 2020, but rather that the database could not collect the vaccination data for that country for 2020.
4. The country records in **Certain\_Country\_Records** can be stored separately to enhance readability (as shown below), as long as the schema format is not changed. However, when there are more certain countries to store in the future, it will be easier to manage if they are stored together in **Certain\_Country\_Records**. The main difference between the two methods is the way queries are written.
  - Wales (iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)
  - Canada (iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)
  - United\_States (iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)
  - Denmark (iso\_code\*, date, total\_vaccinations, people\_vaccinated, people\_fully\_vaccinated, total\_boosters)

5. Due to the lack of detailed information on the number of vaccines provided by each manufacturer in certain regions, In **Provides\_Vaccines\_Certain\_Country**, there is only data on which suppliers provided vaccines daily in that country/location.
6. The original **Vaccinations\_By\_Age\_Group** had many missing values in the age\_group column. To successfully import the CSV file, I removed the rows that lacked data in the age\_group column.