



M1 DEV MANAGER FULLSTACK

# Big Data

"Intégration et Analyse de Données Environnementales :

Une Approche Multi-Dataset pour Comprendre les Défis Globaux"

# Présentation de l'étude sur les données environnementales

## Contexte

Aujourd'hui, les enjeux environnementaux occupent une place centrale dans les débats sociétaux et politiques. L'impact des activités humaines sur l'environnement, notamment en matière d'émissions de CO<sub>2</sub>, de gestion des déchets plastiques et de déforestation, soulève des préoccupations majeures pour la durabilité de notre planète.

Dans ce contexte mondialement préoccupant, analyser et comprendre ces phénomènes à l'aide de données fiables devient indispensable. Ces analyses permettent de mieux cerner les tendances actuelles, de mesurer les impacts environnementaux, et de soutenir la prise de décisions éclairées pour des politiques plus durables.

À travers cette étude, nous avons mobilisé plusieurs ensembles de données riches et variés:

- Les émissions de CO<sub>2</sub> par pays, permettant d'évaluer l'évolution des gaz à effet de serre à l'échelle globale et nationale.
- Les projections et analyses des émissions de CO<sub>2</sub> pour identifier les scénarios futurs possibles et anticiper les impacts environnementaux.
- Les informations sur la surface forestière mondiale, un indicateur clé pour évaluer les pertes en biodiversité et les changements climatiques.

Cette combinaison d'ensembles de données offre une vue d'ensemble des défis environnementaux majeurs et constitue un outil précieux pour explorer des solutions potentielles à ces problématiques complexes.

## Objectif

L'objectif principal de cette étude est de mieux comprendre les impacts environnementaux en analysant des données sur les émissions de CO<sub>2</sub>, la gestion des déchets plastiques et la déforestation. Cela vise à soutenir les initiatives durables et à orienter les décisions politiques et économiques vers des solutions plus respectueuses de l'environnement.

Pour atteindre cet objectif, plusieurs axes d'analyse sont envisagés. Il s'agit tout d'abord de suivre et comparer les émissions de CO<sub>2</sub> par pays afin d'identifier les zones géographiques où des efforts de réduction sont nécessaires. Ensuite, l'évaluation de la gestion des déchets plastiques permettra de mieux comprendre leur ampleur et de proposer des stratégies adaptées aux spécificités locales. L'analyse de l'évolution des surfaces forestières mondiales constitue également un volet clé pour mettre en lumière les régions les plus touchées par la déforestation et ses impacts sur la biodiversité.

Enfin, cette étude proposera des indicateurs de performance environnementale pour mesurer l'efficacité des politiques actuelles. Elle inclura aussi des prédictions basées sur les

## Datasets

- ## Scripts utilisés

## Traitement des données

1. Montage des conteneurs Azure : Les conteneurs **tp-bronze**, **tp-silver**, et **tp-gold** sont montés dans Databricks, permettant d'accéder aux fichiers CSV.

- Chargement des données brutes : Les fichiers CSV du conteneur **tp-bronze** (contenant des données sur les forêts et les émissions) sont lus en Dataframes.

**List csv files**

10:41 PM (20s)

```
mount_name = "tp-bronze"
display(DotUtils.fs.ls(f"/mnt/{mount_name}/emission/"))
display(DotUtils.fs.ls(f"/mnt/{mount_name}/forest_area/"))
```

(8) Spark Jobs

id	path	name	size	modificationTime
1	cdh4/mnt/tp-bronze/emission/Carbon_CO2_Emissions_by_Country.csv	Carbon_CO2_Emissions_by_Country.csv	221702	1732802163000
2	cdh4/mnt/tp-bronze/emission/Emission_By_Country.csv	Emission_By_Country.csv	2957261	1732802163000

2 rows | 12.01 seconds runtime

Refreshed 4 minutes ago

id	path	name	size	modificationTime
1	cdh4/mnt/tp-bronze/forest_area/forest_area_hm.csv	forest_area_hm.csv	68818	1732802181000
2	cdh4/mnt/tp-bronze/forest_area/forest_area_percent.csv	forest_area_percent.csv	99631	1732802181000

2 rows | 12.01 seconds runtime

Refreshed 4 minutes ago

**Get and display all data**

10:41 PM (20s)

```
forest_area_hm = f"/mnt/tp-bronze/forest_area/forest_area_hm.csv"
forest_area_percent = f"/mnt/tp-bronze/forest_area/forest_area_percent.csv"
Carbon_CO2_Emissions_by_Country = f"/mnt/tp-bronze/emission/Carbon_CO2_Emissions_by_Country.csv"
Emission_By_Country = f"/mnt/tp-bronze/emission/Emission_By_Country.csv"

forest_area_hm_df = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load(forest_area_hm)
forest_area_percent_df = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load(forest_area_percent)
Carbon_CO2_Emissions_by_Country_df = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load(Carbon_CO2_Emissions_by_Country)
Emission_By_Country_df = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load(Emission_By_Country)

# display(forest_area_hm_df)
# display(forest_area_percent_df)
# display(Carbon_CO2_Emissions_by_Country_df)
# display(Emission_By_Country_df)
```

(8) Spark Jobs

- Carbon\_CO2\_Emissions\_by\_Country\_df: pyspark.sql.dataframe.DataFrame = [Country: string, Region: string ... 3 more fields]
- Emission\_By\_Country\_df: pyspark.sql.dataframe.DataFrame
  - Country: string
  - ISO 3166-1 alpha-3: string
  - Year: integer
  - Total: double
  - Coal: double
  - Oil: double
  - Gas: double
  - Cement: double
  - Flaring: double
  - Other: double
  - Per Capita: double
- forest\_area\_hm\_df: pyspark.sql.dataframe.DataFrame = [Country Name: string, Country Code: string ... 32 more fields]
- forest\_area\_percent\_df: pyspark.sql.dataframe.DataFrame = [Country Name: string, Country Code: string ... 32 more fields]

- delta\_carbon\_co2\_emissions\_by\_country: pyspark.sql.dataframe.DataFrame
  - country\_name: string
  - region: string
  - kilograms\_of\_co2: double
  - metric\_tons\_per\_capita: double
  - Year: string
- delta\_emission\_by\_country: pyspark.sql.dataframe.DataFrame
  - country\_name: string
  - country\_code: string
  - Year: integer
  - total\_co2\_emission: double
  - coal\_co2\_emission: double
  - oil\_co2\_emission: double
  - gas\_co2\_emission: double
  - cement\_co2\_emission: double
  - flaring\_co2\_emission: double
  - other\_co2\_emission: double
- delta\_forest\_area\_hm: pyspark.sql.dataframe.DataFrame
  - country\_name: string
  - country\_code: string
  - Year: string
  - forest\_area\_hm: double
- delta\_forest\_area\_percent: pyspark.sql.dataframe.DataFrame
  - country\_name: string
  - country\_code: string
  - Year: string
  - forest\_area\_percent: double

### 3. Transformation des données :

- Les DataFrames sont nettoyés, normalisés (renommage des colonnes, suppression des colonnes inutiles, formatage des années), et les données manquantes sont filtrées.
- Les colonnes de type "année" sont pivotées (stack) pour structurer les données de manière cohérente.
- Les jeux de données sont harmonisés avec des noms de colonnes standardisés.

```

Refining the data

# START PM (10)

columns_to_check = ["Coal", "Oil", "Gas", "Coal", "Flaring", "Other", "per_capita"]

delta_emission_by_country = delta_emission_by_country.dropna(subset=columns_to_check, how="all") \
    .filter(delta_emission_by_country["Year"] >= 1990) \
    .drop("per_capita")

# ----- [ FLIPPED YEARS FOR CONSISTENCY ] ----- #
years = [str(year) for year in range(1990, 2022)]
stack_expr = "stack(0), ".format(len(years))

for year in years:
    stack_expr += "'(0)', '(0)', ".format(year)

stack_expr = stack_expr[:-2] + ")" as (Year, forest_area_percent)

delta_forest_area_percent = delta_forest_area_percent.select(
    "country_name",
    "country_code",
    expr(stack_expr)
)

stack_expr = "stack(0), ".format(len(years))

for year in years:
    stack_expr += "'(0)', '(0)', ".format(year)
stack_expr = stack_expr[:-2] + ")" as (Year, forest_area_km)

delta_forest_area_km = delta_forest_area_km.select(
    "country_name",
    "country_code",
    expr(stack_expr)
)

# ----- [ CHANGE YEARS FIELDS FOR CONSISTENCY ] ----- #
delta_carbon_co2_emissions_by_country = delta_carbon_co2_emissions_by_country.withColumn("Year", substring(col("Date"), 7, 4)).drop("Date")

# ----- [ RENAMED FIELDS FOR CONSISTENCY ACROSS DATA ] ----- #

delta_emission_by_country = delta_emission_by_country.withColumnRenamed("TSD_3166-1", "country_code") \
    .withColumnRenamed("Country", "country_name") \
    .withColumnRenamed("Total", "total_co2_emission") \
    .withColumnRenamed("Coal", "coal_co2_emission") \
    .withColumnRenamed("Oil", "oil_co2_emission") \
    .withColumnRenamed("Gas", "gas_co2_emission") \
    .withColumnRenamed("Coal", "coal_co2_emission") \
    .withColumnRenamed("Flaring", "flaring_co2_emission") \
    .withColumnRenamed("per_capita", "per_capita_co2_emission") \
    .withColumnRenamed("Other", "other_co2_emission")

delta_carbon_co2_emissions_by_country = delta_carbon_co2_emissions_by_country.withColumnRenamed("Country", "country_name") \
    .withColumnRenamed("Region", "region")

```

### 4. Stockage des données transformées : Les données sont écrites en format Delta dans le conteneur **tp-silver**, assurant un format optimisé pour des requêtes analytiques.

```

Transform and store initial data to silver delta table

# START PM (10)

mount_name = "tp-silver"
forest_area_km = delta_forest_area_km.withColumnRenamed("Country Name", "country_name").withColumnRenamed("Country Code", "country_code") \
    .write \
    .mode("overwrite") \
    .format("delta") \
    .save("/mnt/{mount_name}/delta_forest_area_km")

forest_area_percent = delta_forest_area_percent.withColumnRenamed("Country Name", "country_name").withColumnRenamed("Country Code", "country_code") \
    .write \
    .mode("overwrite") \
    .format("delta") \
    .save("/mnt/{mount_name}/delta_forest_area_percent")

Carbon_CO2_Emissions_by_Country = delta_carbon_co2_emissions_by_country.withColumnRenamed("kilograms of CO2", "kilograms_of_co2").withColumnRenamed("Metric Tons Per Capita", "metric_tons_per_capita") \
    .write \
    .mode("overwrite") \
    .format("delta") \
    .save("/mnt/{mount_name}/delta_carbon_co2_emissions_by_country")

Emission_By_Country = delta_carbon_co2_emissions_by_country.withColumnRenamed("TSD_3166-1 alpha-1", "TSD_3166-1").withColumnRenamed("Per Capita", "per_capita") \
    .write \
    .mode("overwrite") \
    .format("delta") \
    .save("/mnt/{mount_name}/delta_emission_by_country")

# (S) Spark Jobs

# START PM (11)

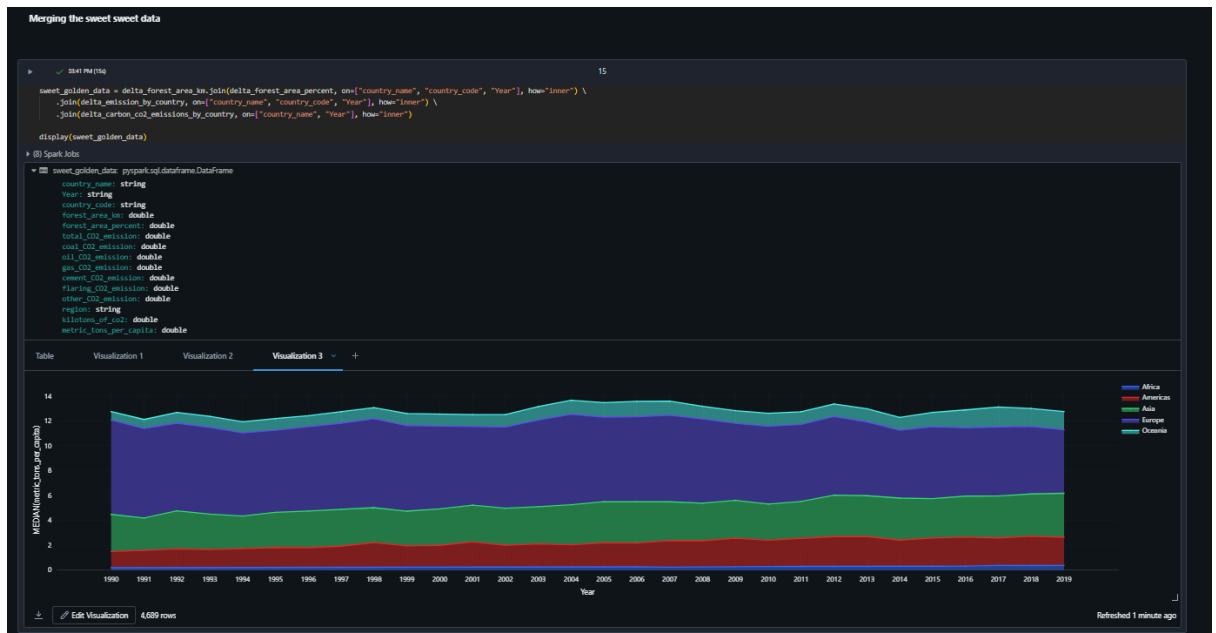
mount_name = "tp-silver"
delta_forest_area_km = spark.read.format("delta").load("/mnt/{mount_name}/delta_forest_area_km")
delta_forest_area_percent = spark.read.format("delta").load("/mnt/{mount_name}/delta_forest_area_percent")
delta_carbon_co2_emissions_by_country = spark.read.format("delta").load("/mnt/{mount_name}/delta_carbon_co2_emissions_by_country")
delta_emission_by_country = spark.read.format("delta").load("/mnt/{mount_name}/delta_emission_by_country")

# display(delta_forest_area_km)
# display(delta_forest_area_percent)
# display(delta_carbon_co2_emissions_by_country)
# display(delta_emission_by_country)

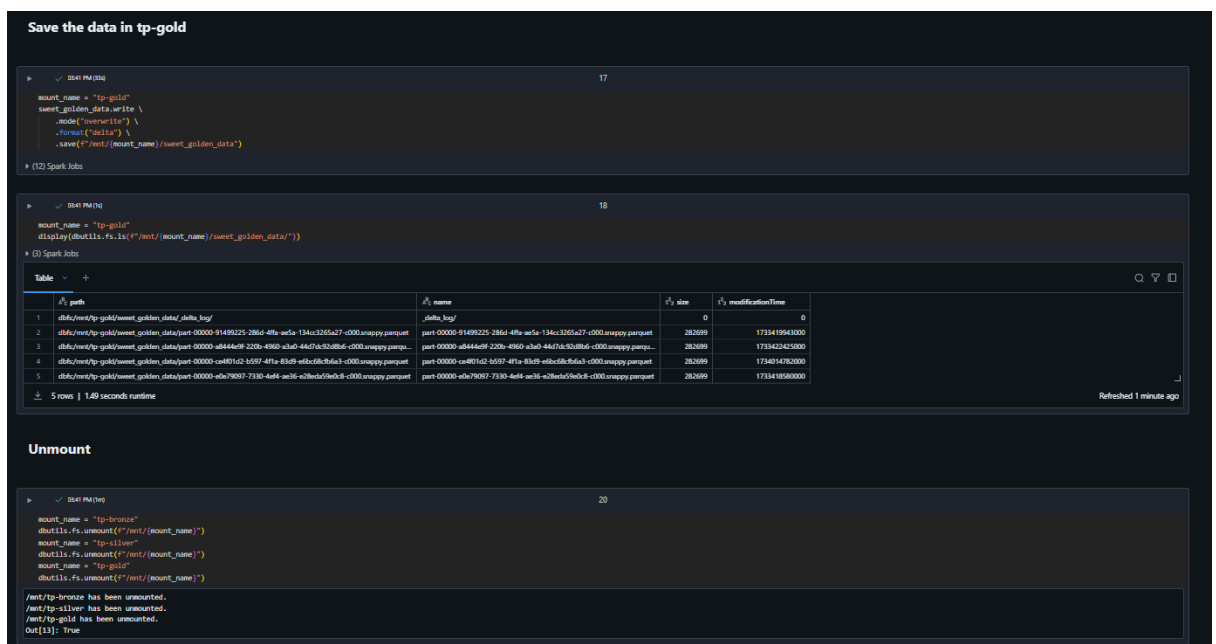
# delta_carbon_co2_emissions_by_country: pyspark.sql.dataframe.DataFrame = [Country: string, Region: string ... 3 more fields]
# delta_emission_by_country: pyspark.sql.dataframe.DataFrame = [Country: string, ISO_3166-1: string ... 9 more fields]
# delta_forest_area_km: pyspark.sql.dataframe.DataFrame = [country_name: string, country_code: string ... 32 more fields]
# delta_forest_area_percent: pyspark.sql.dataframe.DataFrame = [country_name: string, country_code: string ... 32 more fields]

```

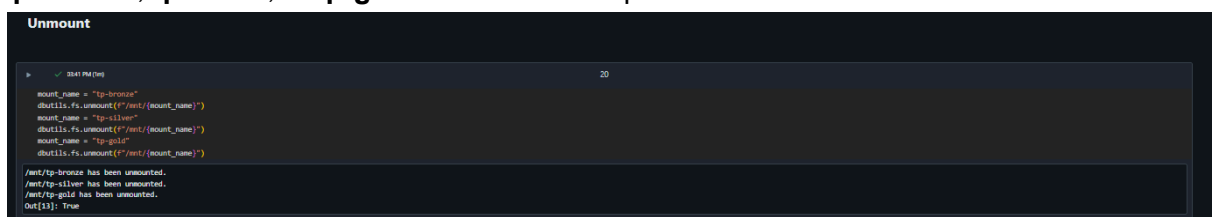
- Fusion des données nettoyées : Un Dataframe unifié, **sweet\_golden\_data**, est créé en joignant les différents jeux de données harmonisés.



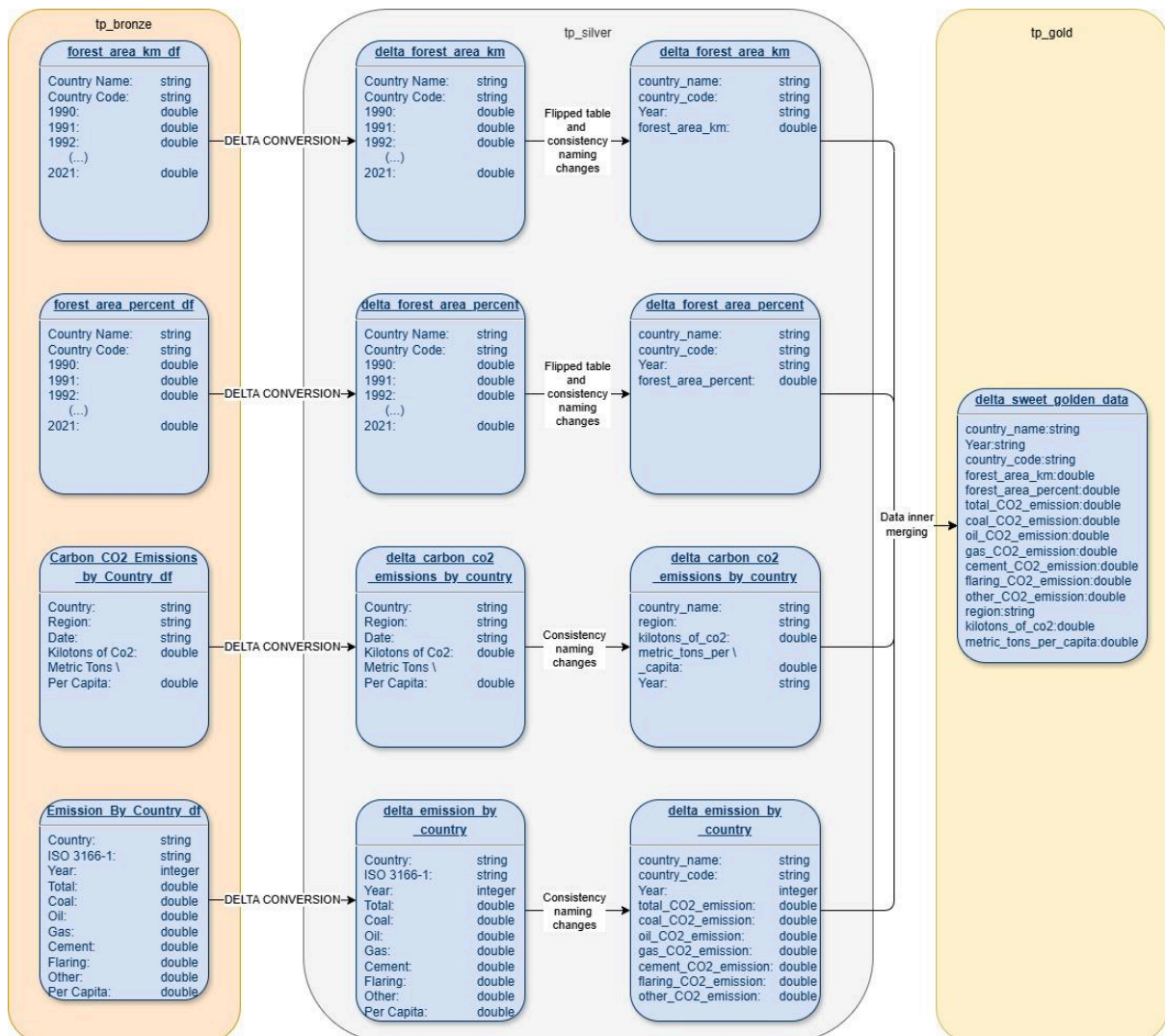
- Stockage final : Le Dataframe fusionné est sauvegardé dans le conteneur **tp-gold** au format Delta.



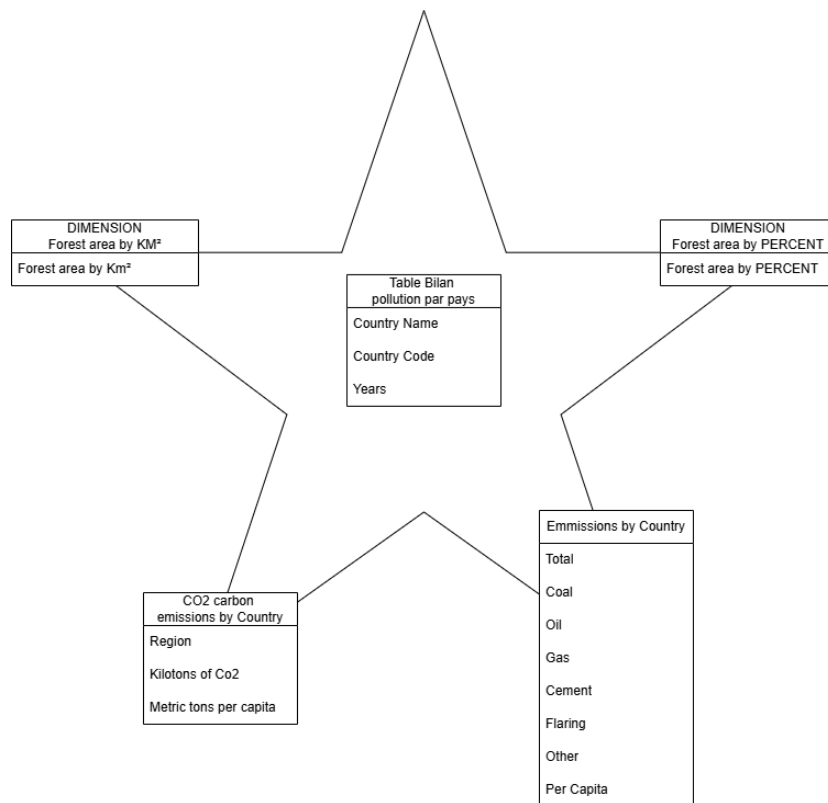
- Démontage des conteneurs : Une fois le traitement terminé, les montages de **tp-bronze**, **tp-silver**, et **tp-gold** sont démontés pour libérer les ressources.



# Modèle Conceptuels des Données



## Modèle en étoile





# Rapport Databricks

screens tableau de bord

