

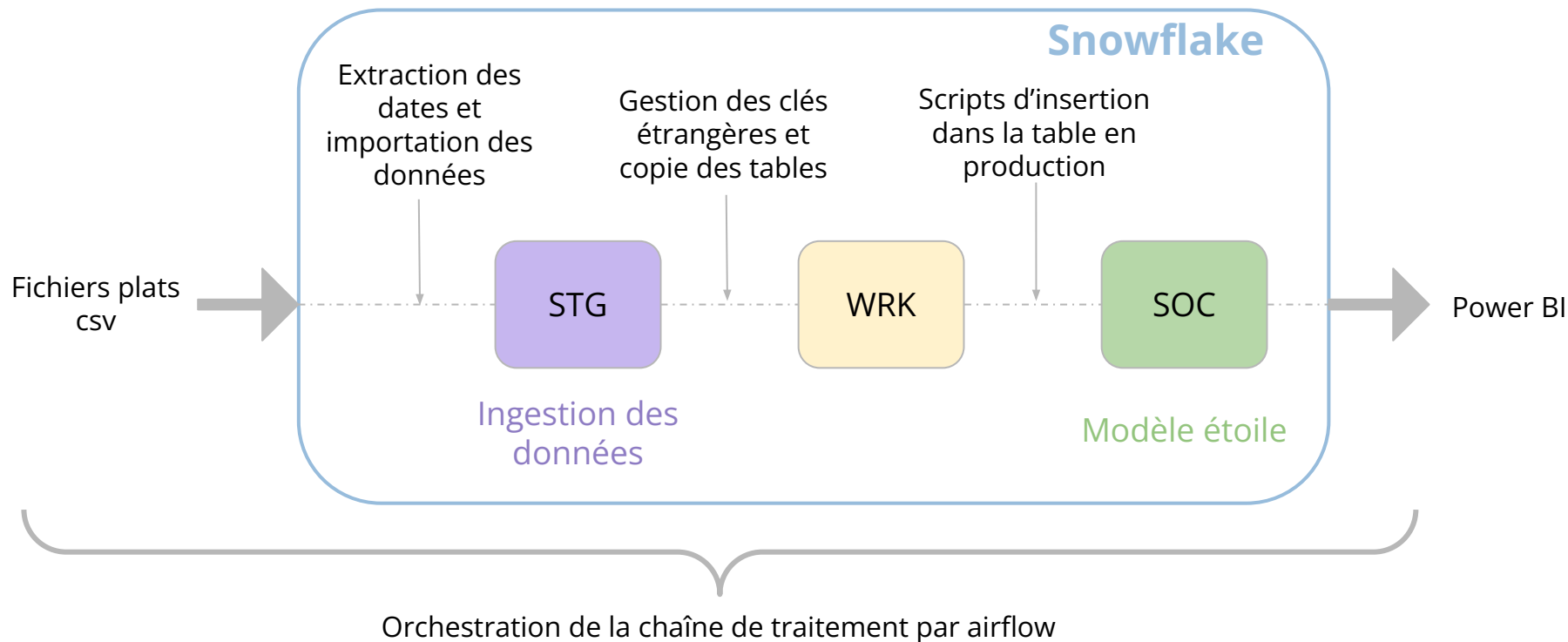
# Projet NF26 en collaboration avec Smart Teem



Mise en place d'une solution  
décisionnelle

# Introduction

## Objectif



# Introduction

Gestion de projet

Organisation 



Développement 



Visual Studio Code



# Présentation des données

## Modèle physique de données du datawarehouse

### Données

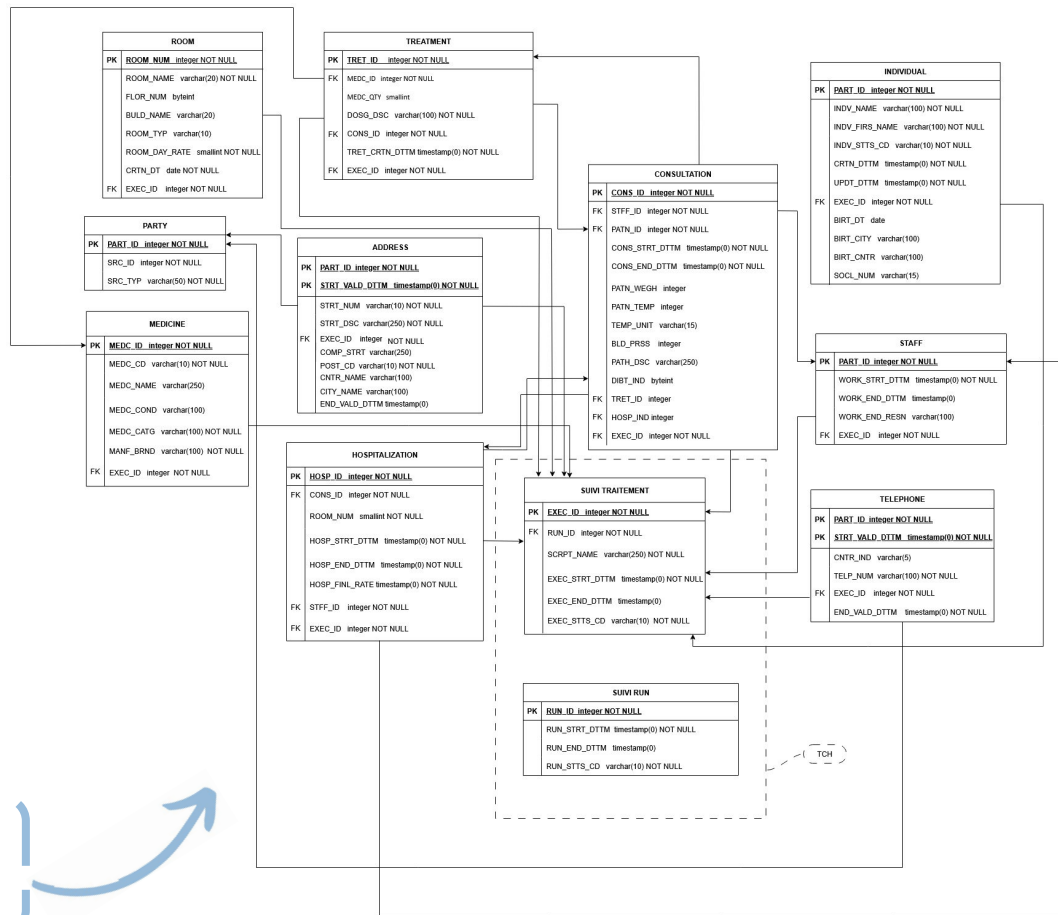
Fichiers .txt journaliers décrivant :

- Chambre,
- Consultation,
- Hospitalisation,
- Médicament,
- Patient,
- Personnel,
- Traitement.

Type de données : nombre, chaîne de caractère, date ...

➤ Date dans un format non reconnu par Snowflake

Création d'un modèle physique basé sur les documents fournis



# Installation du SID

## Création des bases

**Appel du script** create\_database.sql

- Les bases ne sont pas recréées si elles existent déjà

## install\_sid.py

- **Exécution** lancée une seule fois par Airflow
- **Lancement** des scripts SQL de création des tables

## install\_sid.log

- **Exécutions tracées**

## Création des tables

A l'aide de l'excel "Hopital Mapping VF.xlsx" :

1. **Création** des tables dans STG (recréés à chaque exécution)
2. **Création** des tables dans SOC (créés une seule fois)
3. **Création** des tables dans WRK (recréés à chaque exécution)
  - a. SOC + STG

# Ingestion des données

## Outil

### Script de création de procédures

*insert\_STG\_procedure.sql*

- **Instanciation** par airflow pour définir des fonctions d'ingestion simultanée de plusieurs lignes
- **1 procédure** par table de STG
- Permet le **pré-traitement** des données (remplacement des "#" dans les dates par des 0, ajout de valeurs NULL si données d'entrée vide)

Utilise

## Dag pour la pipeline ETL

### Script d'insertion périodique

*launch\_load\_sid.py*

- Permet d'orchestrer la **pipeline de l'ETL** sur un jour donné : **récupération** des données, **transformations** et **ingestion** dans le modèle en production (SOC)
- Vide les tables de **STG** et **WRK** à chaque fois pour éviter les doublons
- Charge les données depuis le csv avec un **script python** puis lance les scripts d'alimentation du datawarehouse avec des **requêtes SQL sur Snowflake**

# Alimentation d'un datawarehouse

## STG → WRK

1. **Sélection** des données de STG
2. **Insertion** dans WRK

## WRK → SOC

1. **Sélection** des données de WRK (SOC)
2. **Suppression** des doublons dans SOC coïncidant (en clé primaire) avec WRK
3. **Insertion** dans SOC

## WRK (STG) → WRK (SOC)

1. **Sélection** des données de WRK (STG)
2. Respect des règles de gestion du fichier Excel
  - a. **Jointures et Restrictions**
  - b. **Surrogate Key** : Séquence incrémentale selon les valeurs uniques de certains attributs
    - i. Window function avec ROW\_NUMBER()
3. **Insertion** des données dans WRK (SOC)



- Permet l'**automatisation journalière** des tâches pour l'ingestion et le traitement des données
- Connexion possible avec **Snowflake**
  - permet l'injection des données en base depuis Airflow
- Utilisation des **DAGs** dans le cadre du projet :
  - un DAG pour la **création** des tables
  - un DAG pour l'**injection** et le **traitement** journaliers des données
- **Gestion des logs** :
  - Automatique
  - Détaillés
  - Pour chaque batch et chaque tâche

## DAG (Directed Acyclic Graph)

Graphe orienté et acyclique utilisé pour modéliser des flux de travail

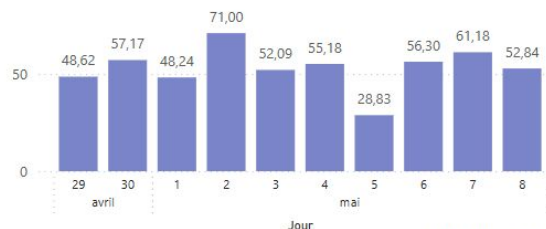


## NF26 - Projet Smart Team

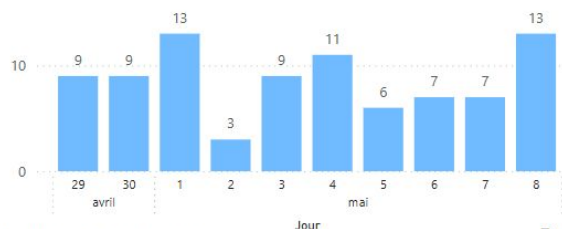
Veuillez sélectionner une pathologie :

Pathology1	Pathology10	Pathology11	Pathology12	Pathology13
Pathology14	Pathology15	Pathology16	Pathology17	Pathology18
Pathology19	Pathology2	Pathology20	Pathology21	Pathology22

Âge moyen des patients



Chambres ayant accueilli des patients diagnostiqués



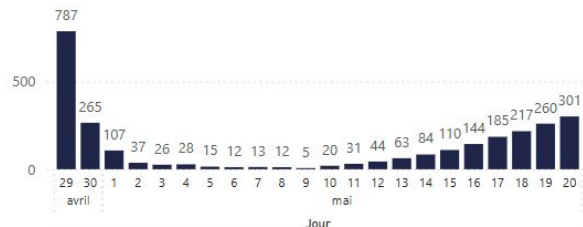
Médicament le plus prescrit

Médicament996

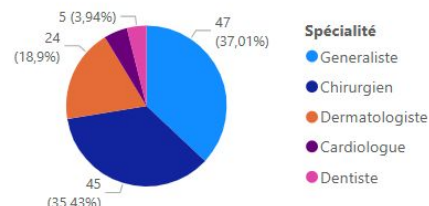
Proportion de patients hospitalisés restés au moins une nuit

100 %

Nombre de chambres inoccupées



Proportion de médecins qui ont diagnostiqué la pathologie



# Power BI

## Requêtes

- Transformations de type
- Jointures
- Colonnes calculées (ex : Âge)
- Colonnes conditionnelles (ex : aPasséNuitHosp)

## Tables ajoutées

- Table de travail (issue des requêtes)
- Calendrier (générée depuis Modèle)
- OccupationJourChambre (générée depuis Modèle)

## Mesures (DAX)

- Chambre inoccupées
- Médicament le plus prescrit
- Proportion de patients hospitalisés restés au moins une nuit

Nom

Table\_de\_travail1

Toutes les propriétés

ÉTAPES APPLIQUÉES

- Source
- O\_INDV (3) développé
- Duplication de la colonne
- Type modifié
- Personnalisée ajoutée
- Type modifié1
- Requêtes fusionnées
- O\_TRET (3) développé
- Requêtes fusionnées1
- R\_MEDC (3) développé
- Requêtes fusionnées2
- O\_HOSP (3) développé
- Requêtes fusionnées3
- R\_PART (3) développé
- Duplication de la colonne1
- Type modifié3
- Duplication de la colonne2
- Type modifié4
- Colonne conditionnelle ajoutée
- Type modifié2

# Conclusion

- Apprentissage d'[Airflow](#) pour l'orchestration des traitements ETL et [Snowflake](#) pour la gestion et le stockage des données
- Gestion de l'entièreté de la chaîne de traitement des données, des fichiers bruts à l'exploitation des données propres avec [Power BI](#)
- Acquisition de compétences techniques en [Data Engineering](#) ainsi que de compétences en travail d'équipe et en collaboration professionnelle avec l'entreprise [Smart Teem](#)