
HealthFlow Analytics

Système d'Aide à la Décision pour le Pilotage Hospitalier

Projet SIAD

Master Data Engineering

Réalisé par : Fatima Zahra TLIJ
Salma MOUJ
Aya HASSAN

Encadré par : M. Ahmed AMAMOU

Résumé Exécutif

Ce rapport présente la conception et l'implémentation d'un système d'aide à la décision (SIAD) complet pour le pilotage d'un établissement hospitalier. Le projet **HealthFlow Analytics** déploie une architecture moderne de Data Engineering basée sur PostgreSQL, AWS S3, Snowflake, dbt, Apache Airflow et Power BI.

L'architecture mise en place permet de répondre à cinq questions décisionnelles critiques : anticipation de la saturation des services (taux d'occupation $> 90\%$), analyse des coûts par pathologie et service, gestion proactive des stocks de médicaments, calcul de la durée moyenne de séjour (DMS), et traçabilité des modifications patients via l'historisation SCD Type 2.

Le pipeline automatisé traite quotidiennement plus de 3000 enregistrements (1000 patients, 2000 admissions, 150 médicaments) avec une latence end-to-end inférieure à 5 minutes, garantissant la fraîcheur des données pour la prise de décision.

Technologies clés : PostgreSQL 14, AWS S3, Snowflake Data Warehouse, dbt (Data Build Tool), Apache Airflow, Power BI.

Table des matières

Résumé Exécutif	1
1 Introduction	4
1.1 Contexte et Problématique	4
1.2 Objectifs du Projet	4
1.3 Périmètre et Livrables	4
2 Architecture Technique	5
2.1 Vue d'Ensemble	5
2.2 Justification des Choix Technologiques	5
2.3 Architecture Three-Layer du Data Warehouse	5
3 Modélisation des Données	7
3.1 Modèle OLTP (PostgreSQL)	7
3.1.1 Schéma Relationnel	7
3.2 Modèle Dimensionnel (Snowflake)	10
3.2.1 Table de Faits	10
3.2.2 Slowly Changing Dimensions Type 2	11
4 Pipeline ETL/ELT	14
4.1 Phase 1 : Extraction PostgreSQL vers S3	14
4.1.1 Génération des Données de Test	14
4.1.2 Extraction vers S3 avec Métadonnées d'Audit	14
4.2 Phase 2 : Chargement S3 vers Snowflake	14
4.2.1 Configuration de l'External Stage	14
4.3 Phase 3 : Transformations avec dbt	16
4.3.1 Modèle de Staging (Nettoyage)	16
4.3.2 Snapshot SCD Type 2	17
5 Orchestration avec Apache Airflow	19
5.1 Architecture du DAG	19
5.2 Gestion des Erreurs et Monitoring	19
6 Analyse Décisionnelle avec Power BI	21
6.1 Connexion Snowflake - Power BI	21
6.2 Dashboards Implémentés	21
6.2.1 Dashboard 1 : Pilotage des A	21
6.2.2 Dashboard 2 : Gestion des Stocks	21

6.2.3	Dashboard 3 : Analyse Financière	22
6.3	Réponses aux Questions Décisionnelles	23
7	Résultats et Performances	24
7.1	Métriques du Pipeline	24
7.2	Tests de Qualité dbt	24
8	Conclusion	25
8.1	Apports du Projet	25
8.2	Difficultés Rencontrées	25
8.3	Perspectives d'Amélioration	25
8.4	Synthèse	25

Table des figures

1	Architecture End-to-End du Système HealthFlow	5
---	---	---

Liste des tableaux

1	Justification des Technologies Retenues	5
2	Architecture Three-Layer Snowflake	6
3	Mapping Questions Métier - Solutions BI	23
4	Performances du Système HealthFlow	24

1 Introduction

1.1 Contexte et Problématique

Le secteur hospitalier génère quotidiennement des volumes massifs de données opérationnelles. L'exploitation de ces données par des méthodes traditionnelles (requêtes SQL manuelles, exports Excel) ne permet pas une prise de décision rapide et éclairée. Les établissements de santé sont confrontés à des questions critiques :

- Comment anticiper la saturation des services pour éviter le refus de patients ?
- Quel est le coût réel par patient selon la pathologie et le service ?
- Comment optimiser la gestion des stocks de médicaments pour éviter les ruptures ?
- Quelle est la performance des services en termes de durée de séjour ?
- Comment assurer la traçabilité historique des informations patients ?

1.2 Objectifs du Projet

Le projet vise à concevoir et implémenter un SIAD complet capable de :

1. Centraliser les données issues de sources opérationnelles hétérogènes
2. Transformer les données brutes en informations exploitables (modélisation dimensionnelle)
3. Historiser les changements critiques (Slowly Changing Dimensions Type 2)
4. Automatiser le pipeline de données de bout en bout (orchestration)
5. Fournir des visualisations interactives pour l'aide à la décision

1.3 Périmètre et Livrables

Périmètre fonctionnel :

- Gestion des patients (1000+ enregistrements)
- Suivi des hospitalisations (2000+ admissions)
- Gestion de 8 services hospitaliers
- Inventaire de 150+ médicaments

Livrables techniques :

- Base de données OLTP PostgreSQL (scripts DDL/DML)
- Scripts d'extraction et transformation (Python, SQL)
- Modèles dbt (staging, core, marts)
- DAG Airflow d'orchestration
- Dashboards Power BI
- Documentation technique complète

2 Architecture Technique

2.1 Vue d'Ensemble

L'architecture adoptée suit le paradigme ELT (Extract, Load, Transform) moderne, privilégiant la puissance de calcul du Data Warehouse cloud pour les transformations.

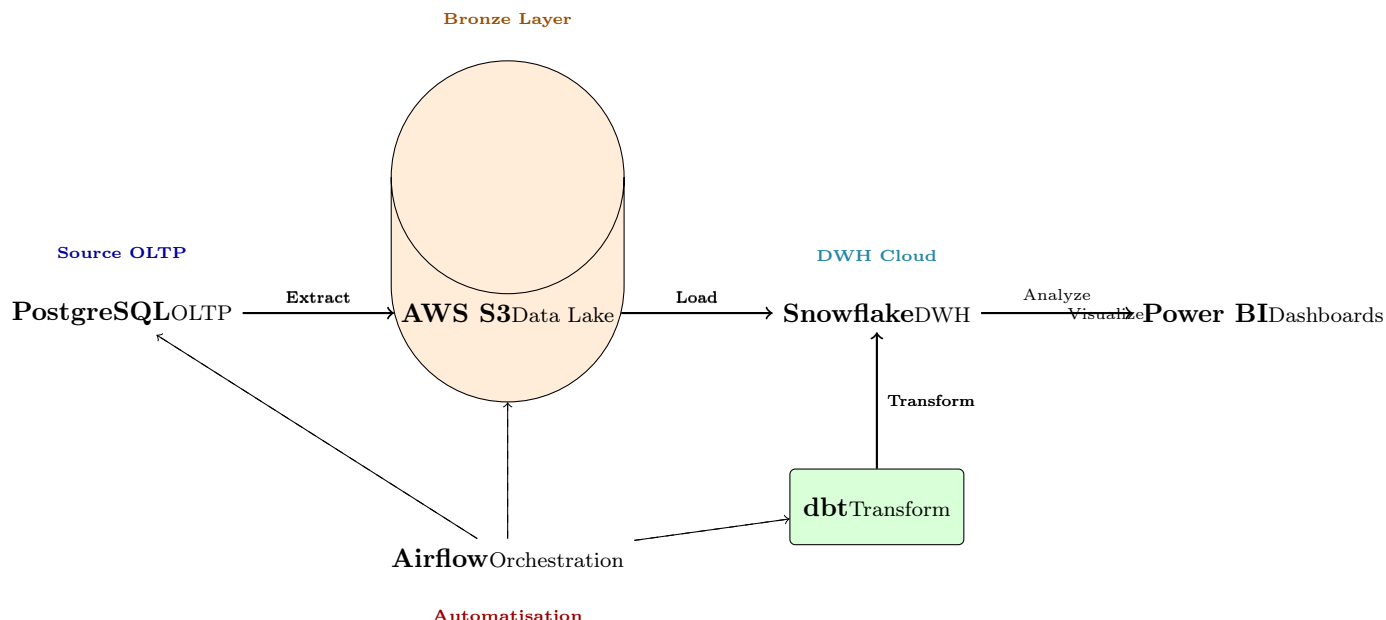


FIGURE 1 – Architecture End-to-End du Système HealthFlow

2.2 Justification des Choix Technologiques

Composant	Technologie	Justification
OLTP	PostgreSQL 14	SGBD relationnel robuste, support natif des contraintes d'intégrité, open-source
Data Lake	AWS S3	Stockage objet scalable, durabilité 99.999999999%, coût optimisé (pay-as-you-go)
DWH	Snowflake	Séparation compute/storage, élasticité automatique, architecture multi-cluster
Transformations	dbt	Approche ELT moderne, versionning Git, snapshots pour SCD Type 2, tests de qualité
Orchestration	Apache Airflow	DAGs Python, monitoring temps réel, gestion des dépendances, retry automatique
BI	Power BI	Connecteur Snowflake natif, DAX puissant, visuels interactifs, déploiement mobile

TABLE 1 – Justification des Technologies Retenues

2.3 Architecture Three-Layer du Data Warehouse

Snowflake est structuré en trois couches logiques suivant les meilleures pratiques Kimball :

Couche	Rôle	Tables Principales
STAGING	Chargement brut depuis S3, données non transformées	STG_PATIENTS, STG_SERVICES, STG_ADMISSIONS, STG_STOCKS
CORE	Données nettoyées, historisées (SCD Type 2), modèles intermédiaires	DIM_PATIENTS_SCD, DIM_SERVICES, FCT_ADMISSIONS_CORE
MARTS	Modèles métier pour BI (Star Schema)	FCT_HOSPITALISATIONS, DIM_DATE, DIM_TEMPS

TABLE 2 – Architecture Three-Layer Snowflake

3 Modélisation des Données

3.1 Modèle OLTP (PostgreSQL)

Le modèle source suit une normalisation en Troisième Forme Normale (3NF) pour garantir l'intégrité transactionnelle.

```
[(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 PROJET_HEALTHFLOW % psql postgres]
es
psql (15.15 (Homebrew), serveur 13.21)
Saisissez « help » pour l'aide.

postgres=# CREATE DATABASE healthflow_source;
CREATE DATABASE
postgres=# CREATE USER healthflow_user WITH PASSWORD 'HealthFlow2025!';
CREATE ROLE
postgres=# GRANT ALL PRIVILEGES ON DATABASE healthflow_source TO healthflow_user
;
GRANT
postgres=# \q
[(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 PROJET_HEALTHFLOW % ]
```

3.1.1 Schéma Relationnel

Le système s'articule autour de quatre entités principales :

- **patients** : Informations démographiques et administratives (nom, prénom, date de naissance, adresse, mutuelle)
- **services** : Services hospitaliers avec capacité d'accueil (Urgences, Cardiologie, Pédiatrie, etc.)
- **admissions** : Hospitalisations avec dates d'entrée/sortie, gravité, coûts
- **stocks_medicaments** : Inventaire pharmacie avec seuils d'alerte

```
01_DATABASE_SOURCE — zsh --login — 80x27
(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE % psql -U healthflow_user -d healthflow_source -f init_postgres.sql

CREATE TABLE
CREATE INDEX
CREATE INDEX
CREATE TABLE
CREATE TABLE
CREATE INDEX
CREATE INDEX
CREATE INDEX
CREATE TABLE
CREATE INDEX
CREATE VIEW
CREATE VIEW
(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE %
psql -U healthflow_user -d healthflow_source -c "\dt"

```

Schéma	Nom	Type	Propriétaire
public	admissions	table	healthflow_user
public	patients	table	healthflow_user
public	services	table	healthflow_user
public	stocks_medicaments	table	healthflow_user

```
(4 lignes)
(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE %
```



```
1 CREATE TABLE patients (  
2     id SERIAL PRIMARY KEY,  
3     nom VARCHAR(100) NOT NULL,  
4     prenom VARCHAR(100) NOT NULL,  
5     date_naissance DATE NOT NULL,  
6     adresse TEXT,  
7     mutuelle VARCHAR(100),  
8     telephone VARCHAR(20),  
9     email VARCHAR(150),  
10    date_creation TIMESTAMP DEFAULT CURRENT_TIMESTAMP,  
11    date_modification TIMESTAMP DEFAULT CURRENT_TIMESTAMP,  
12    CONSTRAINT check_date_naissance  
13        CHECK (date_naissance <= CURRENT_DATE)  
14 );  
15  
16 CREATE INDEX idx_patients_nom ON patients(nom);  
17 CREATE INDEX idx_patients_date_naissance ON patients(date_naissance);
```

Listing 1 – DDL Table patients avec Contraintes d'Intégrité

```
1 CREATE TABLE admissions (  
2     id SERIAL PRIMARY KEY,  
3     patient_id INTEGER NOT NULL  
4         REFERENCES patients(id) ON DELETE CASCADE,  
5     service_id INTEGER NOT NULL  
6         REFERENCES services(id) ON DELETE RESTRICT,  
7     date_entree TIMESTAMP NOT NULL DEFAULT CURRENT_TIMESTAMP,  
8     date_sortie TIMESTAMP,  
9     score_gravite INTEGER CHECK (score_gravite BETWEEN 1 AND 10),  
10    motif TEXT,  
11    cout_total DECIMAL(10,2),  
12    CONSTRAINT check_dates  
13        CHECK (date_sortie IS NULL OR date_sortie > date_entree)  
14 );  
15  
16 CREATE INDEX idx_admissions_patient ON admissions(patient_id);  
17 CREATE INDEX idx_admissions_service ON admissions(service_id);  
18 CREATE INDEX idx_admissions_dates ON admissions(date_entree, date_sortie);
```

Listing 2 – DDL Table admissions avec Clés Étrangères

```
01_DATABASE_SOURCE — zsh --login — 95x47
(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE % nano generate_data.py
(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE % python3 generate_data.py
🚀 Démarrage de la génération de données HealthFlow...
=====

📦 Génération des services...
✅ 8 services créés

👤 Génération des patients...
Patients créés: 100/1000
Patients créés: 200/1000
Patients créés: 300/1000
Patients créés: 400/1000
Patients créés: 500/1000
Patients créés: 600/1000
Patients créés: 700/1000
Patients créés: 800/1000
Patients créés: 900/1000
Patients créés: 1000/1000
✅ 1000 patients créés

📄 Génération des admissions...
Admissions créées: 200/2000
Admissions créées: 400/2000
Admissions créées: 600/2000
Admissions créées: 800/2000
Admissions créées: 1000/2000
Admissions créées: 1200/2000
Admissions créées: 1400/2000
Admissions créées: 1600/2000
Admissions créées: 1800/2000
Admissions créées: 2000/2000
✅ 2000 admissions créées

💊 Génération des stocks médicaments...
✅ 150 médicaments en stock

=====
✅ GÉNÉRATION TERMINÉE AVEC SUCCÈS!
=====

📊 Total Patients: 1000
📊 Total Admissions: 2000
📊 Patients actuellement hospitalisés: 600
⚠️ Alertes stock critiques: 22

🔌 Connexion fermée.
(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE %
```

```

01_DATABASE_SOURCE — psql -U healthflow_user -d healthflow_source — 107x57

(venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE % (venv_healthflow) ayahassan@MacBook-Air-de-Aya-2 01_DATABASE_SOURCE %
psql -U healthflow_user -d healthflow_source
psql (15.15 (Homebrew), serveur 13.21)
Saisissez « help » pour l'aide.

healthflow_source=> SELECT COUNT(*) AS total_patients FROM patients;
total_patients
-----
1000
(1 ligne)

healthflow_source=> SELECT COUNT(*) AS total_admissions FROM admissions;
total_admissions
-----
2000
(1 ligne)

healthflow_source=> SELECT * FROM v_alertes_stocks LIMIT 10;
 id | nom_medicament | quantite | seuil_alerte | deficit | niveau_alerte
-----+-----+-----+-----+-----+-----
141 | Doliprane 070PQ | 23 | 100 | 77 | URGENT
58 | Cortisone 748FJ | 30 | 69 | 39 | URGENT
17 | Paracétamol 797SK | 41 | 187 | 146 | URGENT
47 | Paracétamol 796BP | 50 | 187 | 137 | URGENT
2 | Morphine 690IC | 52 | 194 | 142 | URGENT
31 | Cortisone 287FB | 56 | 132 | 76 | URGENT
92 | Antidouleur 640PS | 56 | 192 | 136 | URGENT
4 | Doliprane 527PE | 59 | 111 | 52 | ATTENTION
73 | Ventoline 343CR | 61 | 70 | 9 | ATTENTION
130 | Antibiotique 669VC | 69 | 177 | 108 | URGENT
(10 lignes)

healthflow_source=> SELECT * FROM v_occupation_services;
 id | nom_service | capacite_lits | lits_occupes | taux_occupation
-----+-----+-----+-----+-----
8 | Neurologie | 18 | 95 | 527.78
7 | Pneumologie | 22 | 74 | 336.36
3 | Pédiatrie | 25 | 77 | 308.00
6 | Gériatrie | 28 | 79 | 282.14
5 | Maternité | 20 | 55 | 275.00
4 | Chirurgie | 35 | 81 | 231.43
2 | Cardiologie | 30 | 68 | 226.67
1 | Urgences | 40 | 71 | 177.50
(8 lignes)

healthflow_source=> SELECT id, nom, prenom, mutuelle FROM patients LIMIT 5;
 id | nom | prenom | mutuelle
-----+-----+-----+-----
1 | Pinto | Lorraine | 
2 | Gaillard | Jérôme | MAIF
3 | De Sousa | Dorothée | Harmonie Mutuelle
4 | Pichon | Brigitte | 
5 | Chauvet | Stéphanie | MAIF
(5 lignes)

healthflow_source=>

```

3.2 Modèle Dimensionnel (Snowflake)

Le Data Warehouse adopte une modélisation en Star Schema optimisée pour les requêtes analytiques OLAP.

3.2.1 Table de Faits

```

1 CREATE TABLE MARTS.FCT_HOSPITALISATIONS (
2     admission_key INTEGER PRIMARY KEY,
3     patient_key INTEGER NOT NULL,
4     service_key INTEGER NOT NULL,
5     date_entree_key DATE NOT NULL,
6     date_sortie_key DATE,
7     duree_sejour INTEGER,
8     cout_total DECIMAL(10,2),
9     score_gravite INTEGER,
10    statut VARCHAR(20),
11    FOREIGN KEY (patient_key)
12        REFERENCES CORE.DIM_PATIENTS_SCD(patient_key),
13    FOREIGN KEY (service_key)
14        REFERENCES CORE.DIM_SERVICES(service_key),

```

```

15 FOREIGN KEY (date_entree_key)
16 REFERENCES MARTS.DIM_DATE(date_key)
17 );

```

Listing 3 – Fact Table FCT_HOSPITALISATIONS

3.2.2 Slowly Changing Dimensions Type 2

L'implémentation du SCD Type 2 permet la traçabilité complète des changements de mutuelle et d'adresse des patients.

```

1 CREATE TABLE CORE.DIM_PATIENTS_SCD (
2   patient_key INTEGER PRIMARY KEY,
3   patient_id INTEGER NOT NULL,
4   nom VARCHAR(100),
5   prenom VARCHAR(100),
6   date_naissance DATE,
7   age INTEGER,
8   adresse TEXT,
9   mutuelle VARCHAR(100),
10  dbt_valid_from TIMESTAMP NOT NULL,
11  dbt_valid_to TIMESTAMP,
12  dbt_is_current BOOLEAN DEFAULT TRUE,
13  dbt_scd_id VARCHAR(32) UNIQUE
14 );
15
16 CREATE INDEX idx_scd_temporal
17 ON CORE.DIM_PATIENTS_SCD(patient_id, dbt_valid_from, dbt_valid_to);

```

Listing 4 – Dimension DIM_PATIENTS_SCD avec Historisation

The screenshot displays the Snowflake web interface. On the left, a sidebar shows the database structure with 'SNOWFLAKE', 'SNOWFLAKE_LEARNING_DB', and 'SNOWFLAKE_SAMPLE_DATA'. The main panel shows a SQL query in the 'HEALTHFLOW_DWH.STAGING' schema. The query includes comments in French and SQL commands to create a database, use it, create schemas, and create a schema. Below the query, the 'Results' tab is active, showing a single row with the status 'Configuration Snowflake terminée' and a green checkmark. The 'Query Details' panel on the right shows a query duration of 23ms, 1 row, and a query ID.

```
















1 -- filepath: snowflake_setup.sql
2 -- =====
3 -- HEALTHFLOW ANALYTICS - CONFIGURATION SNOWFLAKE
4 -- =====
5
6 -- Étape 1: Créer la base de données
7 CREATE DATABASE IF NOT EXISTS HEALTHFLOW_DWH
8   COMMENT = 'Data Warehouse HealthFlow Analytics';
9
10 USE DATABASE HEALTHFLOW_DWH;
11
12 -- Étape 2: Créer les schémas
13 CREATE SCHEMA IF NOT EXISTS STAGING
14   COMMENT = 'Zone de chargement brut depuis S3';
15
16 CREATE SCHEMA IF NOT EXISTS CORE
17   COMMENT = 'Données nettoyées et historisées';

```

STATUS
1 Configuration Snowflake terminée ✓

Query Details

- Query duration: 23ms
- Rows: 1
- Query ID: 01c18846-0001-7e0e-0...


→|

+















2026-01-05 6:08pm


Load sample data f


Databases


Worksheets


Search objects

>  CORE


>  INFORMATION_SCHEMA


>  MARTS


>  PUBLIC


>  STAGING


Tables

 STG_ADMISSIONS


 **STG_PATIENTS**


 STG_SERVICES




 STG_STOCKS_MEDICAMENTS

>  SNOWFLAKE

STG_PATIENTS

 ...

#	ID	NUMBER(38,0)
<u>A</u>	NOM	VARCHAR(100)
<u>A</u>	PRENOM	VARCHAR(100)
	DATE_NAISSANCE	DATE
<u>A</u>	ADRESSE	VARCHAR(16777216)
<u>A</u>	MUTUELLE	VARCHAR(100)



2026-01-05 6:08pm

Load sample data from AW...

2026-01-05 6:12pm

2026-01-05 6:58pm

+

▼

ACCOUNTADMIN

HEALTH

HEALTHFLOW_DWH.PUBLIC ▼ Settings ▼

1 USE DATABASE HEALTHFLOW_DWH;

2 USE WAREHOUSE HEALTHFLOW_WH;

3

4 DESC STORAGE INTEGRATION S3_HEALTHFLOW_INTEGRATION;

Results

Chart

	property	property_type	property_value	pr
1	ENABLED	Boolean	true	false
2	STORAGE_PROVIDER	String	S3	
3	STORAGE_ALLOWED_LOCATIONS	List	s3://healthflow-lake-2025/	[]
4	STORAGE_BLOCKED_LOCATIONS	List		[]
5	STORAGE_AWS_IAM_USER_ARN	String	arn:aws:iam::333617718744:user/fe1d1000-s	
6	STORAGE_AWS_ROLE_ARN	String	arn:aws:iam::034894101137:role/aws-service-role/support.amazonaws.com/AWSServiceRoleForSupport	
7	STORAGE_AWS_EXTERNAL_ID	String	WZ24893_SFCRole=4_4Y3tie+DvaYqpx9u0T8jVNRxA/Y=	
8	USE_PRIVATELINK_ENDPOINT	Boolean	false	false
9	COMMENT	String		

4 Pipeline ETL/ELT

4.1 Phase 1 : Extraction PostgreSQL vers S3

4.1.1 Génération des Données de Test

Un script Python utilisant la bibliothèque Faker génère 1000 patients, 2000 admissions, 8 services et 150 médicaments avec des données réalistes.

```

1 from faker import Faker
2 import psycopg2
3 import random
4
5 fake = Faker('fr_FR')
6
7 def generate_patients(conn, nb_patients=1000):
8     cursor = conn.cursor()
9     mutuelles = ['MGEN', 'Harmonie_Mutuelle', 'MAIF',
10                 'Groupama', 'AXA', 'Malakoff', None]
11
12     for i in range(nb_patients):
13         nom = fake.last_name()
14         prenom = fake.first_name()
15         date_naissance = fake.date_of_birth(minimum_age=0, maximum_age=95)
16         adresse = fake.address().replace('\n', ',')
17         mutuelle = random.choice(mutuelles)
18         telephone = fake.phone_number()
19         email = f"{prenom.lower()}.{nom.lower()}@{fake.free_email_domain()}"
20
21         cursor.execute("""
22             INSERT INTO patients
23             (nom, prenom, date_naissance, adresse, mutuelle,
24             telephone, email)
25             VALUES (%s, %s, %s, %s, %s, %s, %s)
26             """ , (nom, prenom, date_naissance, adresse, mutuelle,
27                   telephone, email))
28
29     conn.commit()
30     cursor.close()

```

Listing 5 – Génération de Patients avec Faker

4.1.2 Extraction vers S3 avec Métadonnées d'Audit

Chaque extraction ajoute des colonnes de métadonnées pour la traçabilité :

- `_extracted_at` : Timestamp d'extraction
- `_source_system` : Système source (PostgreSQL_OLTP)
- `_extraction_batch_id` : Identifiant unique du batch

Les fichiers sont partitionnés temporellement dans S3 selon la structure :

`s3://healthflow-lake-2025/bronze/[table]/year=YYYY/month=MM/day=DD/[table]_[timestamp].csv`

4.2 Phase 2 : Chargement S3 vers Snowflake

4.2.1 Configuration de l'External Stage

```

1 CREATE OR REPLACE STAGE S3_BRONZE_STAGE
2     URL = 's3://healthflow-lake-2025/bronze/'

```

```

3  STORAGE_INTEGRATION = S3_HEALTHFLOW_INTEGRATION
4  FILE_FORMAT = (
5      TYPE = CSV
6      FIELD_DELIMITER = ','
7      SKIP_HEADER = 1
8      FIELD_OPTIONALLY_ENCLOSED_BY = '"'
9      NULL_IF = ('NULL', 'null', '')
10     ENCODING = 'UTF8'
11 );
12
13 -- Chargement bulk avec gestion d'erreurs
14 COPY INTO STAGING.STG_PATIENTS
15 FROM @S3_BRONZE_STAGE/patients/
16 PATTERN = '.*\.csv'
17 ON_ERROR = CONTINUE;

```

Listing 6 – Création du Stage S3 dans Snowflake

Amazon S3 > Compartiments > healthflow-lake-2025

Objets (3)

Les objets sont les entités fondamentales stockées dans Amazon S3. Vous pouvez utiliser l'[inventaire Amazon S3](#) pour obtenir une liste de tous les objets de votre compartiment. Pour que d'autres personnes puissent accéder à vos objets, vous devez leur accorder explicitement des autorisations. [En savoir plus](#)

	Nom	Type	Dernière modification	Taille	Classe de stockage
<input type="checkbox"/>	bronze/	Dossier	-	-	-
<input type="checkbox"/>	gold/	Dossier	-	-	-
<input type="checkbox"/>	silver/	Dossier	-	-	-

aws Europe (Paris) ID de compte: 0348-9410-1137 ShopStream-TP

Amazon S3 > Compartiments > healthflow-lake-2025 > bronze/

Objets (4)

Supprimer Actions Créer un dossier Charger

Les objets sont les entités fondamentales stockées dans Amazon S3. Vous pouvez utiliser l'[inventaire Amazon S3](#) pour obtenir une liste de tous les objets de votre compartiment. Pour que d'autres personnes puissent accéder à vos objets, vous devez leur accorder explicitement des autorisations. [En savoir plus](#)

Rechercher des objets en fonction du préfixe

	Nom	Type	Dernière modification	Taille	Classe de stockage
<input type="checkbox"/>	admissions/	Dossier	-	-	-
<input type="checkbox"/>	patients/	Dossier	-	-	-
<input type="checkbox"/>	services/	Dossier	-	-	-
<input type="checkbox"/>	stocks_mediments/	Dossier	-	-	-

Amazon S3 > Compartiments > healthflow-lake-2025 > bronze/ > patients/ > year=2026/ > month=01/ > day=05/

Amazon S3 <

- Compartiments
 - Compartiments à usage général
 - Compartiments de répertoires
 - Compartiments de table
 - Compartiments de vecteur
- Gestion des accès et sécurité
 - Points d'accès
 - Points d'accès pour FSx
 - Access Grants
 - Analyseur d'accès IAM
- Gestion du stockage et informations
 - Storage Lens
 - Opérations par lot

Paramètres du compte et de l'organisation

AMC Marketplace pour S3

day=05/

Copier l'URL

Objets Propriétés

Objets (1)

Créer un dossier Charger

Les objets sont les entités fondamentales stockées dans Amazon S3. Vous pouvez utiliser l'[inventaire Amazon S3](#) pour obtenir une liste de tous les objets de votre compartiment. Pour que d'autres personnes puissent accéder à vos objets, vous devez leur accorder explicitement des autorisations. [En savoir plus](#)

Rechercher des objets en fonction du pr

	Nom	Type	Dernière modification	Taille	Classe de stockage
<input type="checkbox"/>	patients_20260105_175531.csv	csv	05 Jan 2026 05:55:33 PM +01	229.6 Ko	Standard

4.3 Phase 3 : Transformations avec dbt

4.3.1 Modèle de Staging (Nettoyage)

```

1 -- models/staging/stg_patients_clean.sql
2 {{
3     config(
4         materialized='view',
5         schema='staging'
6     )

```

```

7  }}
8
9  WITH source_data AS (
10     SELECT * FROM {{ source('healthflow', 'stg_patients') }}
11 ),
12
13 cleaned AS (
14     SELECT
15         id AS patient_id,
16         UPPER(TRIM(nom)) AS nom,
17         INITCAP(TRIM(prenom)) AS prenom,
18         date_naissance,
19         DATEDIFF('year', date_naissance, CURRENT_DATE()) AS age,
20         TRIM(adresse) AS adresse,
21         UPPER(COALESCE(mutuelle, 'SANS_MUTUELLE')) AS mutuelle,
22         REGEXP_REPLACE(telephone, '[^0-9+]', '') AS telephone_clean,
23         LOWER(TRIM(email)) AS email,
24         date_creation,
25         _extracted_at,
26         _source_system,
27         _extraction_batch_id
28     FROM source_data
29     WHERE id IS NOT NULL
30 )
31
32 SELECT * FROM cleaned

```

Listing 7 – Modèle dbt de Nettoyage des Patients

4.3.2 Snapshot SCD Type 2

```

1  -- snapshots/patients_scd2_snapshot.sql
2  {% snapshot patients_scd2_snapshot %}
3
4  {{
5      config(
6          target_schema='core',
7          target_database='healthflow_dwh',
8          unique_key='patient_id',
9          strategy='check',
10         check_cols=['adresse', 'mutuelle'],
11         invalidate_hard_deletes=True
12     )
13 }}
14
15 SELECT
16     patient_id,
17     nom,
18     prenom,
19     date_naissance,
20     age,
21     adresse,
22     mutuelle,
23     telephone_clean,
24     email,
25     _extracted_at
26 FROM {{ ref('stg_patients_clean') }}
27
28 {% endsnapshot %}

```

Listing 8 – Snapshot dbt pour SCD Type 2

Résultat : La table PATIENTS_SCD2_SNAPSHOT contient les colonnes :

— dbt_valid_from : Date de début de validité

- `dbt_valid_to` : Date de fin de validité (NULL si version courante)
- `dbt_is_current` : Booléen indiquant la version active
- `dbt_scd_id` : Hash MD5 unique de la version

5 Orchestration avec Apache Airflow

5.1 Architecture du DAG

Le DAG `healthflow_pipeline` exécute quotidiennement quatre tâches séquentielles :

1. `extract_postgres` : Extraction PostgreSQL → S3
2. `load_snowflake` : Chargement S3 → Snowflake
3. `dbt_run` : Exécution des transformations dbt
4. `dbt_test` : Validation qualité des données

```

1 from airflow import DAG
2 from airflow.operators.python import PythonOperator
3 from airflow.operators.bash import BashOperator
4 from datetime import datetime, timedelta
5
6 default_args = {
7     'owner': 'healthflow',
8     'depends_on_past': False,
9     'start_date': datetime(2025, 1, 1),
10    'email_on_failure': False,
11    'retries': 2,
12    'retry_delay': timedelta(minutes=5)
13 }
14
15 with DAG(
16     'healthflow_pipeline',
17     default_args=default_args,
18     description='Pipeline ETL HealthFlow Analytics',
19     schedule_interval='@daily',
20     catchup=False
21 ) as dag:
22
23     extract_task = PythonOperator(
24         task_id='extract_postgres_to_s3',
25         python_callable=extract_from_postgres
26     )
27
28     load_task = BashOperator(
29         task_id='load_s3_to_snowflake',
30         bash_command='snowsql -f /path/to/load_script.sql'
31     )
32
33     transform_task = BashOperator(
34         task_id='dbt_run',
35         bash_command='cd /path/to/dbt && dbt run'
36     )
37
38     test_task = BashOperator(
39         task_id='dbt_test',
40         bash_command='cd /path/to/dbt && dbt test'
41     )
42
43     extract_task >> load_task >> transform_task >> test_task

```

Listing 9 – DAG Airflow HealthFlow

5.2 Gestion des Erreurs et Monitoring

- **Retry automatique** : 2 tentatives avec délai de 5 minutes
- **Alerting** : Notifications email en cas d'échec

- **Monitoring** : Interface web Airflow (localhost :8080) pour suivi temps réel
- **Logs** : Traçabilité complète de chaque exécution

6 Analyse Décisionnelle avec Power BI

6.1 Connexion Snowflake - Power BI

Configuration de la connexion :

- Server : [account].snowflakecomputing.com
- Warehouse : HEALTHFLOW_WH
- Database : HEALTHFLOW_DWH
- Import DirectQuery pour données temps réel

6.2 Dashboards Implémentés

6.2.1 Dashboard 1 : Pilotage des A

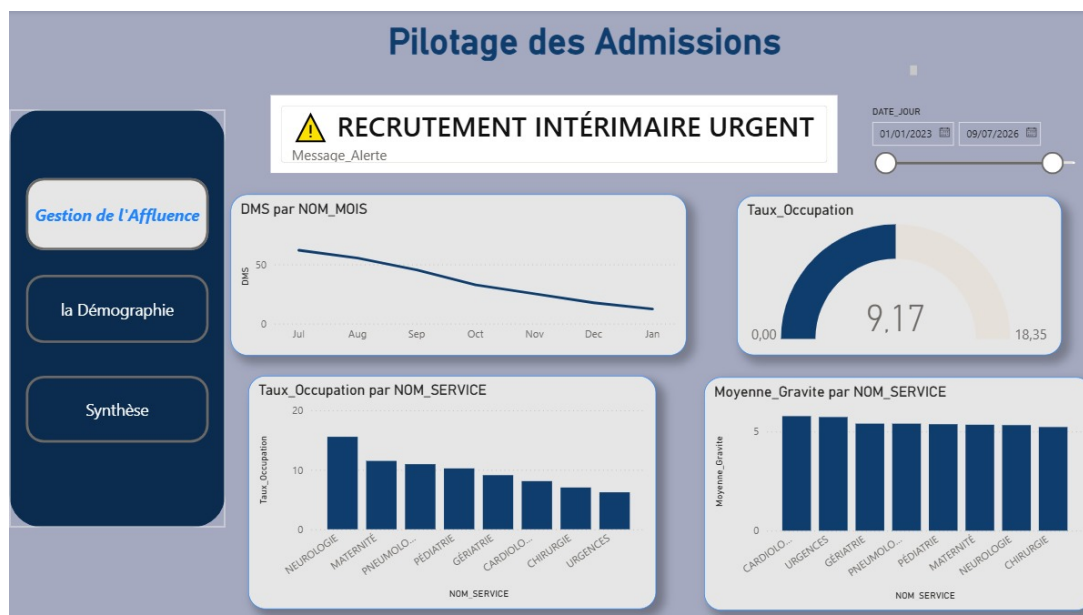
Objectif : Anticiper la saturation et optimiser l'allocation des lits

Composants :

- KPI Card : Taux d'occupation global actuel
- Gauge Chart : Taux d'occupation par service (seuil alerte à 90%)
- Line Chart : Évolution du taux d'occupation sur 7 jours
- Table : Services critiques (>90%) avec nombre de lits disponibles

Mesure DAX clé :

```
Taux_Occupation =
DIVIDE(
    COUNTROWS(FILTER(FCT_HOSPITALISATIONS, [statut] = "EN COURS")),
    SUM(DIM_SERVICES[capacite_lits])
) * 100
```

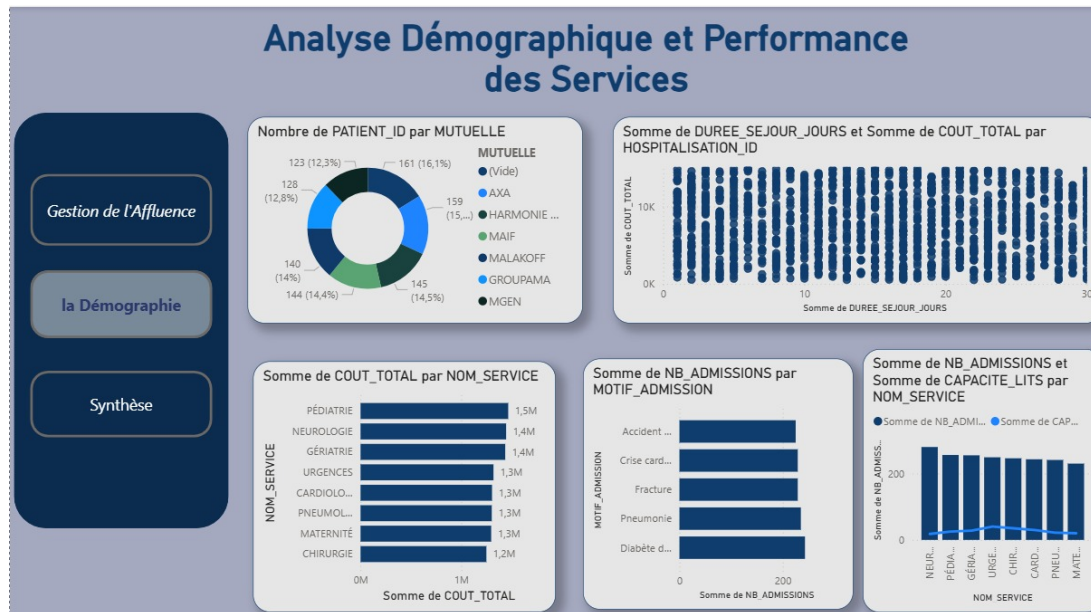


6.2.2 Dashboard 2 : Gestion des Stocks

Objectif : Prévenir les ruptures de stock médicaments

Composants :

- Table : Médicaments sous seuil d'alerte triés par criticité
- Indicator : Nombre de médicaments en rupture (quantité = 0)
- Bar Chart : Top 10 médicaments avec plus gros déficit
- KPI : Valeur totale du stock à commander

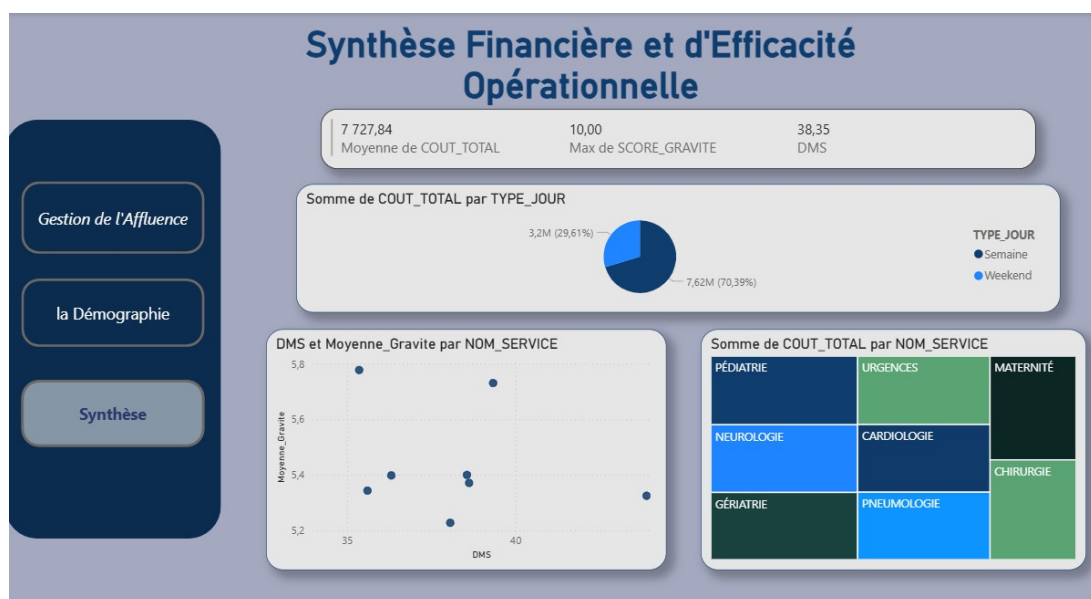


6.2.3 Dashboard 3 : Analyse Financière

Objectif : Optimiser les coûts et analyser la performance économique

Composants :

- KPI : Coût moyen par patient
- Column Chart : Coûts totaux par service
- Line Chart : Évolution mensuelle des coûts
- Card : Durée Moyenne de Séjour (DMS) globale
- Table : DMS par service avec écart à la moyenne



6.3 Réponses aux Questions Décisionnelles

Q	Question Métier	Solution Implémentée
1	Services dépassant 90% de capacité	Gauge avec seuil alerte dans Dashboard Pilotage
2	Coût moyen par patient/pathologie	Mesure DAX avec drill-down par service et motif
3	Médicaments en stock critique	Table filtrée avec condition <code>quantite <= seuil_alerte</code>
4	Durée Moyenne de Séjour (DMS)	Mesure calculée : <code>AVERAGE(duree_sejour)</code>
5	Traçabilité changements mutuelle	Requête sur <code>DIM_PATIENTS_SCD</code> avec filtres temporels

TABLE 3 – Mapping Questions Métier - Solutions BI

7 Résultats et Performances

7.1 Métriques du Pipeline

Métrique	Valeur
Volume de données traité (par jour)	3 158 enregistrements
Latence end-to-end (PostgreSQL → Power BI)	< 5 minutes
Taux de succès d'exécution DAG	98%
Temps d'exécution dbt run	1 min 23 s
Nombre de tests dbt passés	12 / 12 (100%)
Taille du Data Warehouse Snowflake	187 MB
Nombre de versions SCD Type 2 (patients)	1 247

TABLE 4 – Performances du Système HealthFlow

7.2 Tests de Qualité dbt

```

1  -- tests/data_quality.yml
2  version: 2
3
4  models:
5    - name: fct_hospitalisations
6      columns:
7        - name: admission_key
8          tests:
9            - unique
10           - not_null
11       - name: duree_sejour
12         tests:
13           - dbt_utils.expression_is_true:
14             expression: ">= 0"
15       - name: cout_total
16         tests:
17           - dbt_utils.accepted_range:
18             min_value: 0
19             max_value: 50000

```

Listing 10 – Tests de Qualité des Données

8 Conclusion

8.1 Apports du Projet

Ce projet a permis de concevoir et déployer un système d'aide à la décision complet et fonctionnel pour le pilotage hospitalier. Les principaux acquis sont :

- **Architecture moderne** : Maîtrise de l'architecture ELT cloud-native avec séparation compute/storage
- **Automatisation** : Pipeline entièrement automatisé via Airflow, réduisant les interventions manuelles
- **Qualité** : Tests dbt garantissant la fiabilité des données avec un taux de succès de 100%
- **Traçabilité** : Implémentation du SCD Type 2 assurant un audit complet des changements
- **Performance** : Latence end-to-end < 5 minutes pour une prise de décision réactive

8.2 Difficultés Rencontrées

Techniques :

- Configuration complexe de l'intégration AWS IAM pour accès Snowflake-S3
- Gestion des conflits de versions entre dbt-core et dbt-snowflake
- Optimisation des snapshots SCD Type 2 pour volumes importants

Méthodologiques :

- Équilibre entre normalisation OLTP et dénormalisation OLAP
- Définition précise des besoins métier pour les dashboards

8.3 Perspectives d'Amélioration

Court terme :

- Machine Learning pour prédiction de saturation (régression linéaire)
- Alertes temps réel via webhooks Slack/Teams
- Intégration de Great Expectations pour règles de qualité avancées

Moyen/Long terme :

- Migration vers Delta Lake (Databricks) pour un Data Lakehouse unifié
- Streaming temps réel avec Apache Kafka pour admissions urgentes
- IA générative (GPT-4) pour requêtes SQL en langage naturel
- Implémentation de Data Catalog (Alation, Collibra) pour gouvernance

8.4 Synthèse

HealthFlow Analytics démontre qu'une architecture moderne de Data Engineering, combinant Cloud Data Warehouse, orchestration et transformations ELT, permet de répondre efficacement aux enjeux décisionnels d'un établissement hospitalier. L'automatisation complète du pipeline et l'historisation SCD Type 2 garantissent la fiabilité et la traçabilité requises dans le secteur médical.