

# Rapport Professionnel E6

## Evolution d'un Entrepot de Donnees

Projet Data Engineering  
Analyse Territoriale – Region Hauts-de-France

### Competences evaluatees :

- C16** Gerer l'entrepot de donnees a l'aide des outils d'administration et de supervision dans le respect du RGPD, afin de garantir les bons acces, l'integration des evolutions structurelles et son maintien en condition operationnelle dans le temps.
- C17** Implementer des variations dans les dimensions de l'entrepot de donnees en appliquant la methode adaptee en fonction du type de changement demande afin d'historiser les evolutions de l'activite de l'organisation et maintenir ainsi une bonne capacite d'analyse.

**Auteur :** Hamza Elbrek  
**Formation :** Data Engineer  
**Date :** Fevrier 2026

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Methodologie et outillage de maintenance</b>	<b>4</b>
2.1	Gestion des incidents . . . . .	4
2.2	Matrice de priorisation . . . . .	5
<b>3</b>	<b>Journalisation des alertes et des erreurs</b>	<b>5</b>
3.1	Architecture de journalisation mise en place . . . . .	5
3.1.1	Journalisation applicative (Python) . . . . .	5
3.1.2	Journalisation en base de donnees (SQL) . . . . .	5
3.2	Fonction Python de journalisation en base . . . . .	6
3.3	Vue de monitoring des alertes . . . . .	6
3.4	Notifications email automatiques du pipeline ETL . . . . .	7
3.4.1	Configuration . . . . .	7
3.4.2	Email de succes – fin de pipeline . . . . .	7
3.4.3	Email d’alerte – erreur en cours de pipeline . . . . .	8
3.4.4	Statuts des tables dans les emails . . . . .	8
<b>4</b>	<b>Strategie de sauvegarde et restauration</b>	<b>9</b>
4.1	Niveau 1 : Sauvegardes automatiques Azure (PITR) . . . . .	9
4.2	Niveau 2 : Recuperation et export BACPAC vers le Data Lake . . . . .	10
4.2.1	Script d’export BACPAC . . . . .	11
4.2.2	Nettoyage automatique des anciens backups . . . . .	11
4.3	Restauration de la base de donnees . . . . .	11
4.4	Monitoring des sauvegardes . . . . .	11
<b>5</b>	<b>Integration de nouvelles sources de donnees</b>	<b>12</b>
<b>6</b>	<b>Gestion des acces et securite</b>	<b>12</b>
6.1	Roles initiaux du projet E5 . . . . .	12
6.2	Refonte des roles lors de l’evolution E6 . . . . .	13
6.3	Procedure d’integration d’un nouveau collaborateur . . . . .	14
6.4	Row-Level Security (RLS) pour les consultants . . . . .	15
6.4.1	Architecture RLS mise en place . . . . .	15
6.4.2	Schema de securite et tables . . . . .	15
6.4.3	Hierarchie des employes . . . . .	16
6.4.4	Fonction predicat et politique de securite . . . . .	16
6.4.5	Script ETL de chargement du referentiel securite . . . . .	16
<b>7</b>	<b>Procedures d’evolution et de scalabilite</b>	<b>17</b>
<b>8</b>	<b>Variations de dimensions (Slowly Changing Dimensions)</b>	<b>17</b>
8.1	SCD Type 1 : Ecrasement (Overwrite) . . . . .	17
8.2	SCD Type 2 : Historisation (Add New Row) . . . . .	18
8.3	SCD Type 3 : Colonne precedente (Add New Column) . . . . .	19
<b>9</b>	<b>Documentation des evolutions et modeles mis a jour</b>	<b>19</b>
9.1	Evolution des Architectures de Data Warehouse . . . . .	19

9.2	Nouveau modele logique (E6)	20
9.3	Nouveau modele physique (E6)	22
<b>10</b>	<b>Conclusion</b>	<b>23</b>
<b>A</b>	<b>Annexes</b>	<b>24</b>
A.1	Procedures d'exploitation et de scalabilite	24
A.2	Scripts et extraits de code	25
A.3	Glossaire	31
A.4	Arborescence du projet	31

# 1 Introduction

Le projet E5 avait jete les bases d'un entrepot de donnees (Data Warehouse) centralisant les donnees territoriales publiques de la region **Hauts-de-France** : schema en etoile deploye sur **Azure SQL**, six dimensions, sept tables de faits, pipelines ETL Python et infrastructure provisionnee via **Terraform**. Le projet E6 reprend cet entrepot et le fait evoluer en **systeme effectif en conditions reelles** : mise en place d'une methodologie de maintenance, journalisation et monitoring, strategie de sauvegarde, integration de nouvelles sources de donnees, gestion des acces, et implementation des variations de dimensions (Slowly Changing Dimensions) pour historiser les changements au fil du temps.

Pour ancrer le projet dans un contexte metier realiste, un **reseau d'agences bancaires fictif** a ete genere a partir du referentiel geographique de l'entrepot. En effet, les donnees reelles des employes d'une banque constituent des **donnees personnelles sensibles** (RGPD) qu'il est impossible d'importer sur une plateforme cloud Azure dans ce cadre. Les cent une agences ont donc ete positionnees sur les communes de la region dont la population depasse dix mille habitants, et un effectif fictif mais realiste d'environ cinq cent sept employes — directeur regional, directeurs de departement, directeurs d'agence et collaborateurs — a ete genere pour simuler la hierarchie d'un reseau bancaire regional.

Ce referentiel d'employes est le fondement du dispositif de **securite des acces** mis en oeuvre dans E6 : quatre roles orientes metier et une politique Row-Level Security (RLS) sur `dwh.dim_geographie` qui restreint automatiquement la vision des donnees de chaque consultant a son perimetre departemental. Ce rapport documente l'ensemble de ces evolutions et met a jour les modeles logique et physique de l'entrepot.

## 2 Methodologie et outillage de maintenance

### 2.1 Gestion des incidents

Tout incident est initie par la reception d'un message sur la **boite mail generique Outlook** de l'equipe data. Des la reception, un ticket est cree et le processus ci-dessous est applique.

#### 1. Creation du ticket d'incident

Des qu'un email est reçu sur la boite generique de l'equipe, un ticket est ouvert avec les informations suivantes : date de reception, description de la demande, expéditeur, impact present. Aucune demande ne doit rester sans ticket.

#### 2. Prise en charge

A reception du ticket, deux cas sont envisages :

- **Un membre de l'equipe est disponible** → il s'attribue l'incident spontanement et s'en occupe.
- **Aucun membre n'est disponible** → le **manager designe** la personne responsable du traitement.

#### 3. Classification et priorisation (avec le manager)

Le manager et le responsable designe determinent ensemble la priorite de la demande selon la matrice Impact  $\times$  Urgence (cf. section 2.2) : **P1 Critique**, **P2 Haute**, **P3 Moyenne** ou **P4 Basse**. La priorite conditionne le delai de traitement attendu.

#### 4. Accuse de reception vers le metier

La personne en charge contacte le **demandeur ou le referent metier** pour :

- confirmer que la demande est bien prise en compte,
- communiquer la priorite attribuee et une estimation du delai de resolution.

#### 5. Traitement et diagnostic

Investigation du probleme : consultation des logs ETL (`dwh.log_etl`), analyse de la vue `analytics.v_monitoring_alertes`, verification des droits d'accès, identification du composant defaillant. Un journal de bord du diagnostic est maintenu dans le ticket.

#### 6. Resolution

Application du correctif adapte a la nature de l'incident : relance ETL, correction de donnees, ajustement de configuration, mise a jour des droits RLS, etc. La solution est testee et validee avant cloture.

#### 7. Notification de cloture au demandeur

La **personne qui a ouvert l'incident** est informee de la resolution par retour sur le fil mail ou via le ticket. La notification inclut : la cause identifiee, la solution appliquee et la date de cloture.

#### 8. Suivi post-incident et mise en backlog

Si l'incident revele un besoin d'evolution ou un risque recurrent, une tache est ajoutee au **backlog de l'equipe** en respectant la priorisation definie. Les incidents repetitifs font l'objet d'une analyse de cause racine pour eviter les recidives.

*Flux simplifie : Mail reçu → Ticket cree → Prise en charge → Priorite definie → Accuse metier → Diagnostic → Resolution → Notification → Backlog*

## 2.2 Matrice de priorisation

La priorisation des incidents et des demandes de changement repose sur deux criteres : l'**impact** (combien d'utilisateurs ou de processus sont affectes) et l'**urgence** (a quelle vitesse le service doit etre retabli).

Impact / Urgence	Haute	Moyenne	Basse
<b>Eleve</b> (tous les utilisateurs)	P1 – Critique	P2 – Haute	P3 – Moyenne
<b>Moyen</b> (un service/equipe)	P2 – Haute	P3 – Moyenne	P4 – Basse
<b>Faible</b> (un utilisateur)	P3 – Moyenne	P4 – Basse	P4 – Basse

TABLE 1 – Matrice de priorisation Impact  $\times$  Urgence

- **P1 – Critique** : Base de donnees inaccessible, perte de donnees. Resolution immediate ( $< 1h$ ).
- **P2 – Haute** : ETL en echec, datamart non rafraichi. Resolution dans la journee ( $< 8h$ ).
- **P3 – Moyenne** : Performance degradee, acces manquant pour un utilisateur. Resolution sous 48h.
- **P4 – Basse** : Demande d'evolution, amelioration cosmetique. Planification dans le sprint suivant.

## 3 Journalisation des alertes et des erreurs

### 3.1 Architecture de journalisation mise en place

La journalisation est le processus d'enregistrement chronologique des evenements qui se produisent dans un systeme informatique. Notre systeme de journalisation opere a deux niveaux complementaires :

#### 3.1.1 Journalisation applicative (Python)

Les scripts ETL utilisent le module `logging` de Python pour produire des logs structures :

**Code A.1** — *Configuration du logging dans les ETL* ⇒ voir Annexe 13, p. 25

#### 3.1.2 Journalisation en base de donnees (SQL)

En complement des logs applicatifs, une table de journalisation SQL enregistre chaque etape de l'ETL directement dans l'entrepot :

```
1 CREATE TABLE dwh.log_etl (  
2     log_id          INT IDENTITY(1,1) PRIMARY KEY,  
3     date_execution  DATETIME DEFAULT GETDATE(),  
4     etape           NVARCHAR(100),    -- 'load_dimensions', 'load_facts'  
5     ...  
6     table_cible     NVARCHAR(100),    -- 'dwh.dim_temps', 'dwh.'  
7     fait_revenus'
```

```

6      statut          NVARCHAR(20),      -- 'DEBUT', 'SUCCES', 'ERREUR', '
      WARNING'
7      nb_lignes       INT DEFAULT 0,
8      duree_secondes  FLOAT DEFAULT 0,
9      message         NVARCHAR(500),
10     utilisateur     NVARCHAR(100) DEFAULT SYSTEM_USER
11 );

```

Listing 1 – Table de journalisation ETL

**Point important :** La double journalisation (fichier Python + table SQL) assure la traçabilité même si l'un des deux systèmes est indisponible. Les logs Python capturent les erreurs de connexion à la base, tandis que les logs SQL offrent une vue requetable pour les tableaux de bord de supervision.

log_id	date_execution	etape	table_cible	statut	nb_lignes
31	2024-02-11T16:15:54.3100000	BACKUP_BACFAC	DATABASE	SUCCES	0
30	2024-02-11T16:13:20.0170000	BACKUP_BACFAC	DATABASE	DEBUT	0
29	2024-02-11T16:13:20.1770000	load_facts	ALL	SUCCES	0
28	2024-02-11T16:13:20.0500000	load_facts	ALL	DEBUT	0
27	2024-02-11T16:13:20.5800000	load_dimensions	ALL	SUCCES	0
26	2024-02-11T16:13:24.3170000	load_dimensions	ALL	DEBUT	0
25	2024-02-11T16:04:40.2630000	BACKUP_BACFAC	DATABASE	SUCCES	0
24	2024-02-11T16:01:05.5170000	BACKUP_BACFAC	DATABASE	DEBUT	0

FIGURE 1 – Contenu de `dwh.log_etl` – Historique des 20 derniers événements ETL (etape, statut, durée, nombre de lignes)

## 3.2 Fonction Python de journalisation en base

**Code A.2** — Fonction `log_etl_db()` dans les scripts ETL

⇒ voir Annexe 14, p. 26

## 3.3 Vue de monitoring des alertes

Une vue analytique permet de superviser les alertes et erreurs depuis un outil BI :

```

1 CREATE VIEW analytics.v_monitoring_alertes AS
2 SELECT
3     CAST(date_execution AS DATE) AS jour,
4     etape,
5     statut,
6     COUNT(*) AS nb_evenements,
7     SUM(CASE WHEN statut = 'ERREUR' THEN 1 ELSE 0 END) AS nb_erreurs,
8     AVG(duree_secondes) AS duree_moyenne_sec
9 FROM dwh.log_etl
10 GROUP BY CAST(date_execution AS DATE), etape, statut;

```

## Listing 2 – Vue de monitoring pour la supervision

date	table	status	metric1	metric2	metric3	metric4	metric5
2026-02-20	load_dimensions	SKIP	18	0	0	0.025761666666666666	0
2026-02-20	load_dimensions	SUCCESS	3	0	0	0.05459423333333333	0
2026-02-20	load_facts	DEBUT	3	0	0	0	0
2026-02-20	load_facts	ERREUR	3	3	0	0.38119833333333333	0
2026-02-20	load_facts	ERREUR_PARTIELLE	2	0	0	2.6874335	0
2026-02-20	load_facts	IGNORER	1	0	0	0	0
2026-02-20	load_facts	SKIP	14	0	0	0.33805407142857137	0
2026-02-20	load_facts	SUCCESS	1	0	0	3.365682	0
2026-02-22	load_dimensions	DEBUT	1	0	0	0	0
2026-02-22	load_dimensions	SKIP	6	0	0	0.0234755	0

FIGURE 2 – Contenu de `dw.h.log_erreurs` – Historique des 20 dernières erreurs détectées durant les exécutions ETL

### 3.4 Notifications email automatiques du pipeline ETL

En complément de la journalisation en base, un module de notification email (`etl_notifier.py`) envoie automatiquement un rapport à la fin de chaque exécution du pipeline ETL, qu'elle soit réussie ou en erreur.

#### 3.4.1 Configuration

La configuration SMTP est stockée dans `terraform.tfvars` pour centraliser tous les paramètres du projet. Les variables d'environnement ont priorité si elles sont définies :

```
1 # E6 - Notifications email ETL
2 etl_smtp_host      = "smtp.gmail.com"
3 etl_smtp_port      = "587"
4 etl_smtp_user      = "moncompte@gmail.com"
5 etl_smtp_password  = "xxxx xxxx xxxx xxxx" # App Password Gmail
6 etl_notify_email   = "destinataire@example.com"
```

Listing 3 – Configuration SMTP dans `terraform.tfvars`

**Point important :** Pour Gmail, il faut générer un **App Password** dans Sécurité → Mots de passe des applications (après activation de la vérification en 2 étapes). Ce mot de passe à 16 caractères remplace le mot de passe Google dans la configuration ETL.

#### 3.4.2 Email de succès – fin de pipeline

À la fin d'un pipeline complet sans erreur, l'orchestrateur (`run_etl.py`) appelle `send_success_email` qui envoie un rapport HTML structuré :



```

1 # En cas de succes global du pipeline
2 if success:
3     rapport_etapes['details'] = rapport_details
4     send_success_email(rapport_etapes, smtp_config=smtp_config)

```

Listing 4 – Declenchement de la notification dans run\_etl.py

L'email de succes contient deux tableaux :

- **Recapitulatif des etapes** : statut (OK/ERREUR), nombre de lignes et duree pour chacune des 5 etapes (Staging, Dimensions, Faits, Refresh, Backup)
- **Detail par table** : statut (OK, SKIP, IGNORE, ERREUR), heure et duree pour chaque table chargee (dimensions + faits)

### 3.4.3 Email d'alerte – erreur en cours de pipeline

Des qu'une erreur est detectee sur une table ou une etape, l'orchestrateur envoie immediatement un email d'alerte sans attendre la fin du pipeline :

```

1 send_error_email(
2     etape='Faits',
3     table='dwh.fait_revenus',
4     erreur="KeyError: colonne 'obs_value' introuvable",
5     heure='14:32:18',
6     rapport_partiel={ # Tables deja traitees avant l'erreur
7         'dwh.fait_population': {'statut': 'OK', 'nb_lignes': 45},
8         'dwh.fait_entreprises': {'statut': 'SKIP', 'nb_lignes': 0},
9     },
10    smtp_config=smtp_config,
11 )

```

Listing 5 – Envoi d'alerte en cas d'erreur ETL (run\_etl.py)

L'email d'alerte contient :

- Un bloc rouge avec : l'etape, la table en erreur, l'heure et le message d'erreur complet
- Un tableau des tables deja traitees avant l'erreur (avec leurs statuts)

### 3.4.4 Statuts des tables dans les emails

Statut	Icone	Signification
OK	✓	Table chargee avec succes
SKIP	→	Table ignoree (deja a jour, pas de nouvelles donnees)
IGNORE	→	Table ignoree (table staging source manquante)
ERREUR	×	Erreur lors du chargement

TABLE 2 – Statuts des tables dans les emails de notification ETL

## Explication technique : Architecture de notification

Le module `etl_notifier.py` utilise uniquement la bibliotheque standard Python (`smtplib`, `email.mime`) – aucune dependance externe necessaire. La connexion SMTP est etablie en **TLS** (port 587) avec **STARTTLS**. Si la configuration SMTP est absente ou incomplete, les notifications sont simplement ignorees (log **WARNING**) sans bloquer le pipeline.

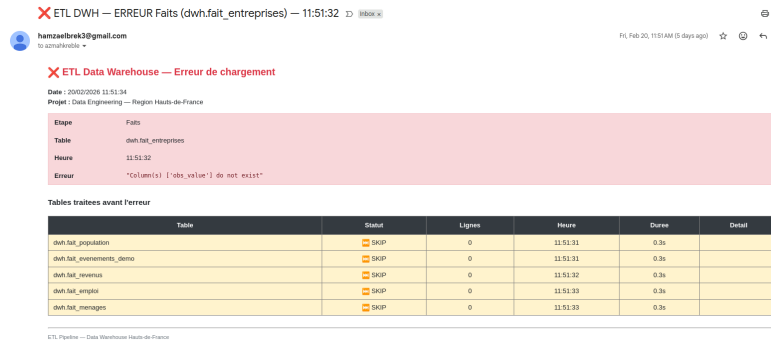


FIGURE 3 – Exemple de notification email automatique – Rapport de fin de pipeline ETL (succes ou erreur)

## 4 Strategie de sauvegarde et restauration

### 4.1 Niveau 1 : Sauvegardes automatiques Azure (PITR)

Azure SQL Database effectue automatiquement des sauvegardes, sans aucune intervention de notre part. Ce mecanisme, appele **PITR** (Point-In-Time Restore), constitue le premier niveau de notre strategie.

La politique de retention est definie dans notre infrastructure Terraform :

```
1 resource "azurerm_mssql_database" "sql" {
2   # ...
3   short_term_retention_policy {
4     retention_days          = 14      # PITR sur 14 jours
5     backup_interval_in_hours = 12    # Differentiel toutes les 12h
6   }
7
8   long_term_retention_policy {
9     weekly_retention       = "P4W"   # 4 semaines
10    monthly_retention      = "P12M"  # 12 mois
11    yearly_retention       = "P3Y"   # 3 ans
12    week_of_year           = 1
13  }
14 }
```

Listing 6 – Configuration Terraform de la retention des sauvegardes (main.tf)

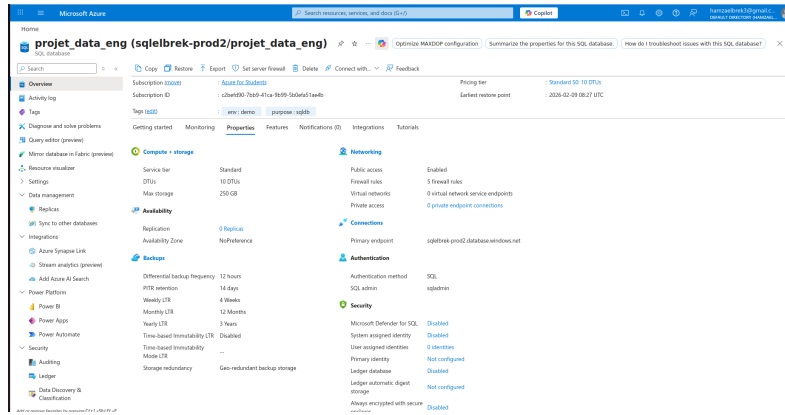


FIGURE 4 – Configuration des sauvegardes dans le portail Azure – Politique PITR (14 jours) et retention long terme (4 semaines, 12 mois, 3 ans)

## 4.2 Niveau 2 : Recuperation et export BACPAC vers le Data Lake

Les sauvegardes automatiques Azure sont gérées par la plateforme et restent dans l'écosystème Azure. Pour disposer d'une copie **indépendante** et **portable** de notre base, nous ajoutons un second niveau : après chaque exécution de l'ETL, un script Python récupère l'état de la base et l'exporte en fichier `.bacpac` horodaté vers notre Data Lake ADLS Gen2.

Le principe est le suivant :

1. Azure réalise automatiquement ses backups (complet + différentiel + logs de transactions)
2. Notre script `backup_to_datalake.py`, intégré comme **étape 5 du pipeline ETL**, exporte la base via `az sql db export` vers le Data Lake
3. Le fichier `.bacpac` obtenu est stocké dans le container `raw/backups/` avec un horodatage
4. Ce fichier peut être utilisé à tout moment pour restaurer la base via `az sql db import`

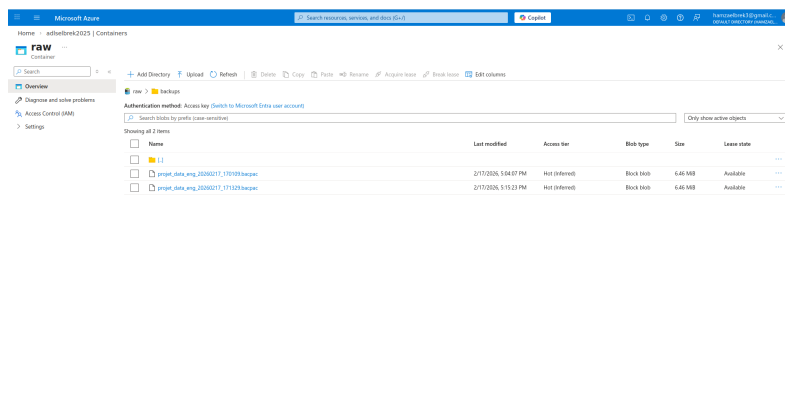


FIGURE 5 – Fichiers `.bacpac` stockés dans le Data Lake – Container `raw/backups/` avec horodatage

### 4.2.1 Script d'export BACPAC

**Code A.4** — Export BACPAC vers le Data Lake (*backup\_to\_datalake.py*) ⇒ voir Annexe 16, p. 26

### 4.2.2 Nettoyage automatique des anciens backups

Le script gere egalement la suppression des fichiers `.bacpac` de plus de 30 jours (configurable) pour eviter l'accumulation dans le Data Lake :

**Code A.5** — Nettoyage automatique des anciens backups ⇒ voir Annexe 17, p. 27

## 4.3 Restauration de la base de donnees

En cas de besoin (corruption, perte de donnees, incident), deux methodes de restauration sont disponibles selon la situation :

- **Restauration Point-in-Time (PITR)** : Utilise les sauvegardes automatiques Azure pour restaurer la base a n'importe quel instant dans la fenetre de retention (14 jours). Methode la plus rapide, ideale pour les incidents recents.
- **Import BACPAC depuis le Data Lake** : Recree la base a partir d'un fichier `.bacpac` exporte par notre pipeline ETL. Utile pour restaurer un etat au-dela de la fenetre PITR, ou pour migrer la base vers un autre serveur.

**Code A.6** — Restauration via Azure CLI

⇒ voir Annexe 18, p. 27

## 4.4 Monitoring des sauvegardes

Deux vues SQL permettent de surveiller l'etat des sauvegardes :

- `analytics.v_etat_backup_azure` : Interroge la DMV `sys.dm_database_backups` pour afficher les sauvegardes automatiques Azure (type, date, taille).
- `analytics.v_historique_backups` : Filtre la table `dwh.log_etl` pour afficher l'historique des exports BACPAC effectues par le pipeline ETL.

```

1 CREATE VIEW analytics.v_historique_backups AS
2 SELECT
3     log_id, date_execution,
4     etape AS type_backup, statut, message,
5     duree_secondes, nb_lignes, utilisateur
6 FROM dwh.log_etl
7 WHERE etape IN ('BACKUP_BACPAC', 'RESTAURATION');
```

Listing 7 – Vue de l'historique des exports BACPAC

## 5 Integration de nouvelles sources de donnees

L'entrepot E5 ne chargeait que 4 tables de faits sur les 7 prevues. Pour E6, deux nouvelles sources CSV ont ete integrees : `EMPLOI_CHOMAGE_hauts_de_france.csv` alimentant `fait_emploi`, et `Menage_hauts_de_france.csv` alimentant `fait_menages`. Le processus ETL suit un pattern uniforme : les donnees sont d'abord chargees depuis la zone de staging, les clees etrangeres sont resolues via les dimensions `dim_temps` et `dim_geographie`, les valeurs sont agregees par annee et departement, puis inserees dans la table de faits. Chaque chargement est trace dans le journal `log_etl` en base de donnees.

emploi_id	temps_id	geo_id	dems_id	population_active	population_en_emploi
1	1	1	13	481678.137079999	488239.01086
2	1	2	13	2318460.23682	1899752.87426
3	1	3	13	773238.88888	685450.81261
4	1	4	13	1273472.83552	107698.00954
5	1	5	13	518217.50356	485752.86427

FIGURE 6 – Donnees chargees dans `dwh.fait_emploi` – Exemple des 5 premieres lignes apres integration de la nouvelle source

## 6 Gestion des acces et securite

La gestion des acces constitue un pilier fondamental de la gouvernance de l'entrepot de donnees. Cette section decrit l'evolution du dispositif de securite entre le projet initial (E5) et la version actuelle (E6), ainsi que les procedures operationnelles associees.

### 6.1 Roles initiaux du projet E5

Dans le projet initial E5, l'entrepot de donnees disposait de trois roles RBAC (*Role-Based Access Control*) couvrant les grands profils d'utilisation :

Role	Permissions	Utilisateurs types
role_etl_process	Lecture/ecriture sur <code>stg</code> et <code>dwh</code> , lecture sur <code>dm</code> et <code>analytics</code>	Comptes de service ETL
role_analyst	Lecture complete sur <code>dwh</code> (dimensions et faits), <code>dm</code> et <code>analytics</code> . Aucun acces au staging.	Data Analysts
role_dwh_admin	Controle total sur tous les schemas et objets	Administrateurs DBA

TABLE 3 – Roles RBAC du projet initial E5

Ce dispositif répondait correctement aux besoins d'un entrepot centralise avec des profils generiques. Cependant, lors du passage en production avec les 519 collaborateurs repartis sur les 5 departements de la region Hauts-de-France, une **limite structurelle** est apparue : ces roles ne prenaient pas en compte les restrictions geographiques.

En effet, un collaborateur de l'agence de Valenciennes (Nord, dept. 59) disposant du role `role_analyst` avait techniquement acces aux donnees de toute la region — y compris les departements de l'Oise (60) et de la Somme (80) dans lesquels il n'a aucune responsabilite. Cette situation est incompatible avec le principe RGPD de *minimisation des acces*.

## 6.2 Refonte des roles lors de l'evolution E6

Pour repondre a cette limite, l'evolution E6 a redesigne le dispositif de securite selon deux axes complementaires :

1. **Refonte des roles RBAC** : suppression du role `role_dwh_admin`, scission du role `role_analyst` en deux profils distincts, et ajout d'un role `role_admin` avec permissions explicites.
2. **Introduction du RLS** (*Row-Level Security*) : filtrage automatique des donnees par zone geographique pour les collaborateurs en agence.

Role	Permissions	Utilisateurs types
<code>role_admin</code>	Contrôle total (CONTROL) sur tous les schemas. Peut créer des tables, vues, procédures.	Administrateurs DBA
<code>role_etl_process</code>	Lecture/écriture (SELECT, INSERT, UPDATE, DELETE) sur <code>stg</code> et <code>dwh</code> . Lecture seule sur <code>dm</code> et <code>analytics</code> .	Comptes de service ETL
<code>role_analyst</code>	Lecture seule sur <code>dwh</code> (dimensions + faits), <code>dm</code> et <code>analytics</code> . Aucun accès au staging.	Data Analysts internes
<code>role_consultant</code>	Lecture seule sur <code>dm</code> et <code>analytics</code> uniquement. <b>Accès filtre par département via RLS</b> (voir section 6.4).	Collaborateurs agences

TABLE 4 – Roles RBAC refondus lors de l'evolution E6

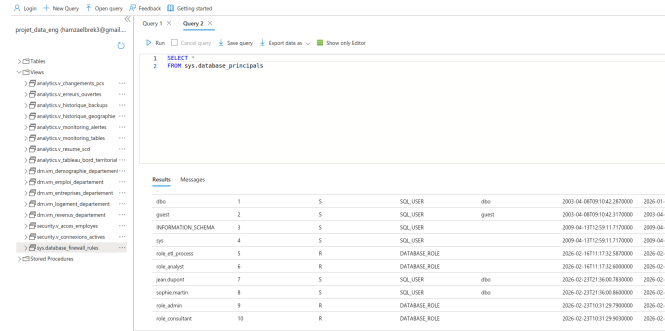


FIGURE 7 – Roles de base de donnees crees dans Azure SQL – Resultat de SELECT name FROM sys.database\_principals WHERE type = 'R'

### 6.3 Procedure d'integration d'un nouveau collaborateur

Lorsqu'un nouveau collaborateur rejoint la structure, l'integration de ses acces a l'entrepot suit un processus en trois etapes selon son profil.

#### Etape 1 — Determiner le role approprié

Profil	Role a attribuer
Administrateur / DBA	role_admin
Ingenieur ETL / pipeline de donnees	role_etl_process
Data Analyst (acces region entiere)	role_analyst
Collaborateur en agence (acces departement)	role_consultant

TABLE 5 – Correspondance profil metier / role SQL

#### Etape 2 — Creer l'utilisateur et attribuer le role

```

1 -- Creer l'utilisateur avec mot de passe (Azure SQL Database)
2 CREATE USER [prenom.nom] WITH PASSWORD = 'MotDePasseSecurise!';
3
4 -- Attribuer le role correspondant au profil
5 ALTER ROLE role_consultant ADD MEMBER [prenom.nom];

```

Listing 8 – Creation d'un nouvel acces utilisateur dans Azure SQL

#### Etape 3 — Configurer la zone géographique (si role\_consultant)

Pour un collaborateur en agence, il faut enregistrer son perimetre géographique dans la table security.utilisateurs\_zones afin que le RLS filtre correctement ses acces :

```

1 -- Acces restreint a un departement (ex : Nord - 59)
2 INSERT INTO security.utilisateurs_zones (login_sql, departement_code)
3 VALUES ('prenom.nom', '59');
4
5 -- Acces region entiere (ex : directeur regional)
6 INSERT INTO security.utilisateurs_zones (login_sql, departement_code)
7 VALUES ('prenom.nom', NULL);

```

Listing 9 – Enregistrement de la zone géographique du collaborateur

## Explication technique : Logique du filtre RLS

La fonction predicat `security.fn_rls_geographie` verifie la presence du login dans `utilisateurs_zones`. Si le collaborateur n'y figure pas, il a acces a toutes les lignes (cas des roles `role_admin`, `role_analyst`, `role_etl_process`). S'il y figure avec un `departement_code` non null, il ne voit que les donnees de ce departement. Avec NULL, il voit toute la region.

## 6.4 Row-Level Security (RLS) pour les consultants

### 6.4.1 Architecture RLS mise en place

Le RLS s'applique sur la table `dwh.dim_geographie` (qui contient le `departement_code`). Puisque toutes les vues des schemas `dm` et `analytics` effectuent une jointure sur cette table, le filtre se propage automatiquement a toutes ces vues sans modification supplémentaire.

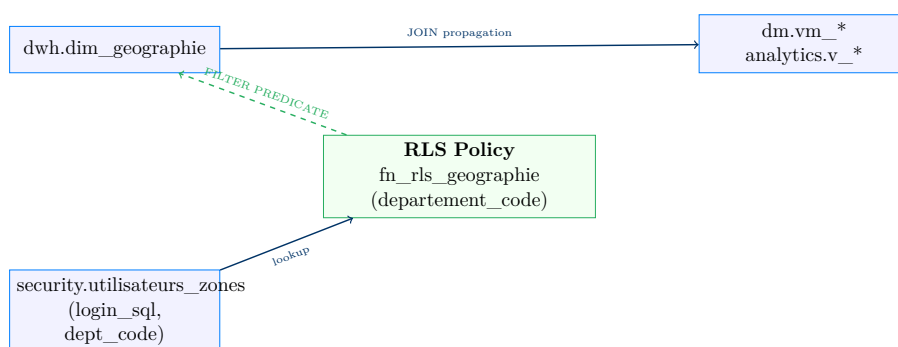


FIGURE 8 – Architecture RLS – Filtre sur `dim_geographie`, propagation aux datamarts

Le screenshot montre l'interface Azure Data Studio avec une requête SQL exécutée dans l'éditeur de requêtes. La requête est `SELECT * FROM [dm].[dim_geographie] WHERE [departement_code] = 59`. Les résultats de la requête sont affichés dans le volet des résultats.

annee	departement_code	departement_nom	population_totale	nbhabitants	dmom	nbde_natur
2010	59	Nord	10461305.8611726			
2015	59	Nord	10462307.91396132	3641192	2387700	1253490
2021	59	Nord	10470622.984739998	3211370	2648855	501215

FIGURE 9 – Verification du filtre RLS en base – Un consultant du departement 59 ne voit que les donnees de son perimetre

### 6.4.2 Schema de securite et tables

Un schema `security` dedie contient les trois tables du dispositif RLS :



### 6.4.3 Hierarchie des employes

La structure hierarchique compte 4 niveaux pour les 5 departements de la region (02 Aisne, 59 Nord, 60 Oise, 62 Pas-de-Calais, 80 Somme) :

Niveau	Effectif	Perimetre RLS
DIRECTEUR_REGIONAL	1	Region entiere (5 departements)
DIRECTEUR_DEPARTEMENT	5	Leur departement
DIRECTEUR_AGENCE	101	Departement de leur agence
COLLABORATEUR	412	Departement de leur agence
<b>Total</b>		<b>519 employes</b>

TABLE 6 – Hierarchie des employes et perimetre d'accès RLS

Les 101 agences correspondent aux communes de la region dont la population depasse 10 000 habitants (9 grandes agences >50 000 hab., 41 moyennes, 51 petites). Le script Python `analytics/etl/load_security.py` genere automatiquement cette structure a partir du referentiel geographique `communes.json`.

### 6.4.4 Fonction predicat et politique de securite

#### Explication technique : Propagation automatique du filtre

Le predicat RLS est applique sur `dwh.dim_geographie`. Toutes les vues `dm.*` et `analytics.*` joignent cette table via `geo_id`. Lorsqu'un consultant execute une requete sur `dm.vm_demographie_departement`, SQL Server filtre d'abord les lignes de `dim_geographie` selon son predicat, puis effectue la jointure avec les tables de faits. Le resultat est automatiquement restreint a ses departements autorises, sans modifier les vues ni les requetes.

### 6.4.5 Script ETL de chargement du referentiel securite

Le script `analytics/etl/load_security.py` permet de generer et charger les donnees de securite :

```

1 # Verification sans chargement en base
2 python load_security.py --check
3
4 # Chargement initial (ou rechargement complet)
5 python load_security.py --reset --load
6
7 # Resultat :
```

```

8 # 101 agences -> security.agences
9 # 519 employes -> security.employes
10 # 519 zones -> security.utilisateurs_zones

```

Listing 10 – Utilisation de load\_security.py

## 7 Procédures d’évolution et de scalabilité

L’entrepôt de données est conçu pour évoluer sans interruption de service. Quatre procédures opérationnelles encadrent les opérations courantes d’administration :

- **Procédure 1** : Créer un nouvel accès utilisateur
- **Procédure 2** : Ajouter un nouveau datamart analytique
- **Procédure 3** : Montée en charge (scalabilité)
- **Procédure 4** : Ajouter une nouvelle source de données

Le détail complet de chaque procédure est documenté en annexe :

**Annexe B** — *Procédures d’exploitation et de scalabilité* ⇒ voir Annexe A.1, p. 24

## 8 Variations de dimensions (Slowly Changing Dimensions)

### 8.1 SCD Type 1 : Ecrasement (Overwrite)

#### Definition : SCD Type 1

Le Type 1 consiste à **écraser purement et simplement** l’ancienne valeur par la nouvelle. L’historique est perdu.

**Exemple** : Le libelle de la section NAF “C” est corrigé de “Industrie manufacturere” (faute) à “Industrie manufacturiere”.

**Code A.10** — *Implementation SCD Type 1 – Ecrasement* ⇒ voir Annexe 21, p. 29

## 8.2 SCD Type 2 : Historisation (Add New Row)

### Definition : SCD Type 2

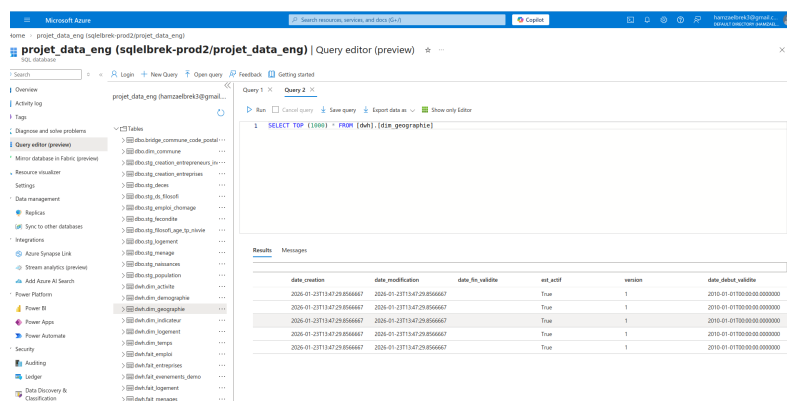
Le Type 2 est l'approche la plus riche : au lieu d'écraser l'ancienne valeur, on **creé une nouvelle ligne** pour la nouvelle valeur et on **ferme l'ancienne ligne** en marquant sa date de fin. Chaque version de la dimension coexiste dans la table avec des colonnes de gestion :

- `date_debut_validite` : Date a partir de laquelle cette version est valide
- `date_fin_validite` : Date de fin de validite (NULL = version courante)
- `est_actif` : Booleen indiquant si c'est la version en vigueur
- `version` : Numero de version incremental

**Exemple** : La commune “Noyelles-Godault” (code 62617) fusionne avec “Henin-Beaumont” en 2026. On conserve l'historique des deux noms.

```
1 -- Ajout des colonnes SCD Type 2 a dim_geographie
2 ALTER TABLE dwh.dim_geographie ADD
3     date_debut_validite DATETIME DEFAULT '2010-01-01',
4     date_fin_validite   DATETIME NULL, -- NULL = actif
5     est_actif          BIT DEFAULT 1,
6     version            INT DEFAULT 1;
```

Listing 11 – Colonnes SCD Type 2 ajoutes a dim\_geographie



date_creation	date_modification	date_fin_validite	est_actif	version	date_debut_validite
2026-01-23T13:47:29.856667	2026-01-23T13:47:29.856667		True	1	2010-01-01T00:00:00.000000
2026-01-23T13:47:29.856667	2026-01-23T13:47:29.856667		True	1	2010-01-01T00:00:00.000000
2026-01-23T13:47:29.856667	2026-01-23T13:47:29.856667		True	1	2010-01-01T00:00:00.000000

FIGURE 10 – Colonnes SCD Type 2 dans dwh.dim\_geographie – Champs `est_actif`, `version` et `date_debut_validite` pour l'historisation des communes

**Code A.11** — *Procédure SCD Type 2 – Historisation*

⇒ voir Annexe 22, p. 30

## 8.3 SCD Type 3 : Colonne precedente (Add New Column)

### Definition : SCD Type 3

Le Type 3 est un compromis entre le Type 1 et le Type 2 : on ajoute une **colonne supplémentaire** pour stocker l'ancienne valeur, tout en mettant à jour la valeur courante. Cette approche ne conserve qu'un **seul niveau d'historique** (la valeur precedente).

- `libelle_actuel` : La valeur courante de l'attribut
- `libelle_precedent` : L'ancienne valeur (avant le dernier changement)
- `date_changement` : La date du dernier changement

**Exemple** : La categorie PCS “Cadres” est renommee en “Cadres et professions intellectuelles superieures”.

```
1  -- Ajout des colonnes SCD Type 3 a dim_demographie
2  ALTER TABLE dwh.dim_demographie ADD
3      ancien_libelle_pcs NVARCHAR(200) NULL,
4      date_changement_pcs DATETIME NULL;
```

Listing 12 – Colonnes SCD Type 3 ajoutees a dim\_demographie

**Code A.12** — *Procedure SCD Type 3*

⇒ voir Annexe 23, p. 30

## 9 Documentation des evolutions et modeles mis a jour

### 9.1 Evolution des Architectures de Data Warehouse

Le modele logique E6 s'articule desormais autour de **cinq schemas** avec des responsabilites clairement separees :



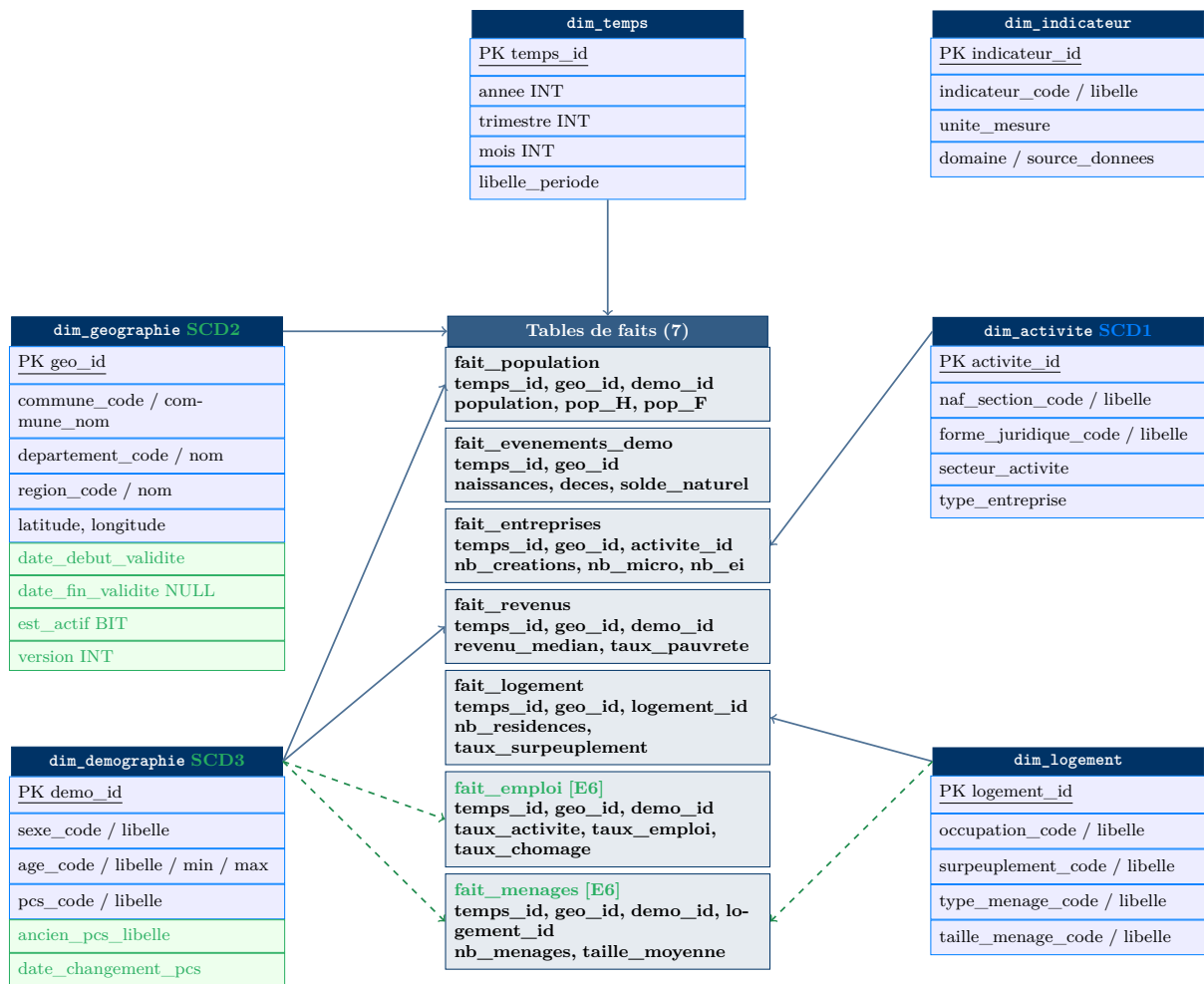


FIGURE 13 – Schema en etoile complet – Data Warehouse E6 (vert = nouveautes E6, tirete = nouvelles liaisons)

### 9.3 Nouveau modele physique (E6)

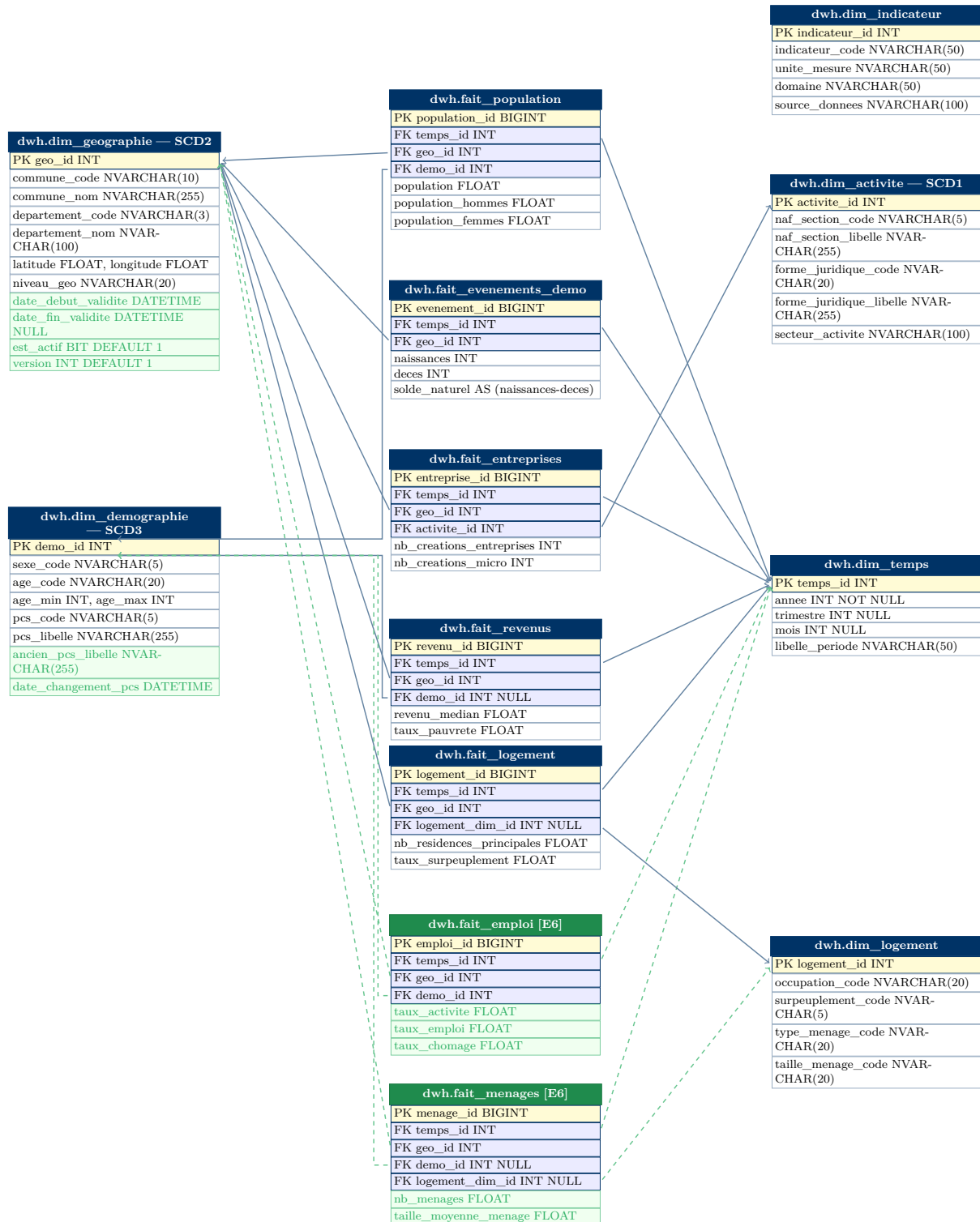


FIGURE 14 – Modele physique ERD – Data Warehouse E6 complet (jaune = PK, bleu clair = FK, vert = nouveautes E6, fleche pleine = lien E5, fleche tiretee = nouveau lien E6)

## 10 Conclusion

Le projet E6 a transformé l'entrepôt de données E5 en un système opérationnel complet. Sur le plan de la maintenance, un processus de gestion des incidents calqué sur la réalité de l'équipe a été formalisé, accompagné d'une journalisation à deux niveaux (logs Python et table `dwh.log_etl`) et de notifications email automatiques à l'issue de chaque exécution du pipeline. Une stratégie de sauvegarde avec export BACPAC vers le Data Lake ADLS Gen2 et des politiques de rétention Azure complètent ce dispositif. L'entrepôt a également été étendu par l'intégration de deux nouvelles tables de faits (emploi, ménages) et la refonte des rôles de sécurité : quatre rôles RBAC métier et un mécanisme de Row-Level Security sur `dwh.dim_geographie` qui restreint automatiquement les données visibles par chaque consultant à son périmètre départemental.

Sur le plan dimensionnel, les trois types de variations de Kimball (SCD 1, 2 et 3) ont été implémentés sur des dimensions distinctes, avec les procédures SQL `MERGE` et les adaptations ETL correspondantes. La simulation d'un effectif fictif de cinq cent sept employés répartis dans cent une agences ancrées dans la géographie régionale a permis de valider concrètement le dispositif de sécurité dans un contexte métier réaliste. L'ensemble des évolutions est documenté dans ce rapport avec la mise à jour des modèles logique et physique de l'entrepôt.

Pour aller plus loin, plusieurs évolutions pourraient prolonger ce travail : intégration d'Azure Monitor pour des alertes temps réel, mise en place d'un Change Data Capture (CDC) pour déclencher les ETL à la source, adoption d'un framework de qualité des données (Great Expectations, dbt), ou migration du pipeline vers Azure Data Factory pour une orchestration native et un monitoring centralisé.



## A Annexes

### A.1 Procédures d'exploitation et de scalabilité

#### Procédure 1 – Créer un nouvel accès utilisateur

1. Recevoir la demande d'accès via le système de ticketing (type : `changement`, priorité : P4)
2. Identifier le rôle approprié selon le profil : `role_analyst` (analyste interne) ou `role_consultant` (collaborateur agence avec RLS)
3. Vérifier que l'employé existe dans `security.employees`; si absent, l'ajouter avec le script `load_security.py --load`
4. Créer l'utilisateur SQL et assigner le rôle via le script ETL : `python load_security.py --create-users`
5. Vérifier les permissions effectives avec `security.v_acces_employees` et tester la connexion depuis Azure Data Studio
6. Documenter dans le ticket : login créé, rôle attribué, zone RLS, date
7. Notifier l'utilisateur de ses identifiants par canal sécurisé

#### Procédure 2 – Ajouter un nouveau datamart analytique

1. Recevoir le besoin d'analyse (ticket type : `changement`)
2. Identifier les tables de faits et dimensions impliquées
3. Créer la vue dans le schéma `dm` en jointurant sur les dimensions actives (`est_actif = 1` si SCD Type 2 implique)
4. Accorder les droits de lecture : `GRANT SELECT ON dm.<vue> TO role_consultant`
5. Valider la propagation du filtre RLS si la vue jointure sur `dim_geographie`
6. Ajouter un test dans `test_e6_evolution.py` pour vérifier l'existence de la vue
7. Documenter la vue dans le dictionnaire de données et mettre à jour le modèle logique

**Code A.9** — *Exemple de création de vue datamart*

⇒ voir Annexe A.2, p. 29

#### Procédure 3 – Montée en charge (scalabilité)

1. **Diagnostic** : Identifier le goulot d'étranglement via les vues de monitoring (`analytics.v_resume_et_sys.dm_exec_sessions`)
2. **Base de données** : Augmenter le SKU Azure SQL dans `terraform.tfvars` (ex : S0 → S2 → S4) puis `terraform apply`
3. **Partitionnement** : Si les tables de faits dépassent 10 millions de lignes, envisager un partitionnement par année sur `annee_id`
4. **Index** : Ajouter des index couvrants sur les colonnes de filtrage fréquentes (`departement_code`, `annee_id`)
5. **ETL** : Activer le chargement parallèle par département dans `run_etl.py` pour réduire la durée du pipeline

6. **Archivage** : Exporter les donnees anterieures a N-3 vers le Data Lake via `backup_to_datalake.py` et les supprimer du DWH actif
7. **Validation** : Contrôler les performances apres chaque modification et documenter les gains dans le ticket de changement

#### Procedure 4 – Ajouter une nouvelle source de donnees

1. **Analyse** : Etudier le format (CSV, API, JSON), le volume et la frequence de mise a jour de la nouvelle source
2. **Staging** : Creer la table de staging `stg.stg_nouvelle_source` avec les colonnes brutes de la source
3. **Mapping** : Identifier les correspondances avec les dimensions existantes (ou creer de nouvelles dimensions si necessaire)
4. **Schema** : Creer ou modifier la table de faits dans le schema `dwh` avec les cles etrangeres vers les dimensions
5. **ETL** : Developper la fonction de chargement dans `load_facts.py` en respectant le pattern SCD adequat pour chaque dimension impliquee
6. **Datamart** : Creer ou mettre a jour les vues dans le schema `dm` (cf. Procedure 2)
7. **Tests** : Ajouter les tests de validation dans `test_e6_evolution.py`
8. **Deploiement** : Integrer le script SQL dans le pipeline Terraform (`Terraform/sql/`) et valider avec `terraform apply`
9. **Documentation** : Mettre a jour le modele logique et le dictionnaire de donnees

## A.2 Scripts et extraits de code

### A.1 — Configuration du logging dans les ETL

```
1 import logging
2
3 # Configuration du logger avec sortie fichier + console
4 logging.basicConfig(
5     level=logging.INFO,
6     format='%(asctime)s [%(levelname)s] %(name)s - %(message)s',
7     handlers=[
8         logging.FileHandler('etl_pipeline.log'),
9         logging.StreamHandler() # Sortie console
10    ]
11 )
12 logger = logging.getLogger('etl_pipeline')
13
14 # Exemples d'utilisation
15 logger.info("Debut du chargement de dim_geographie")
16 logger.warning("Table staging vide, etape ignoree")
17 logger.error("Cle etrangere invalide pour fait_revenus")
```

Listing 13 – Configuration du logging dans les ETL

## A.2 — Fonction log\_etl\_db() dans les scripts ETL

```
1 def log_etl_db(engine, etape, table_cible, statut,
2               nb_lignes=0, duree=0, message=''):
3     """Enregistre un evenement ETL dans dwh.log_etl."""
4     try:
5         with engine.connect() as conn:
6             conn.execute(text("""
7                 INSERT INTO dwh.log_etl
8                     (etape, table_cible, statut, nb_lignes,
9                     duree_secondes, message)
10                    VALUES (:etape, :table, :statut, :nb,
11                            :duree, :msg)
12                    """), {
13                        'etape': etape, 'table': table_cible,
14                        'statut': statut, 'nb': nb_lignes,
15                        'duree': duree, 'msg': message
16                    })
17             conn.commit()
18     except Exception as e:
19         logger.error(f"Impossible de loguer en BDD: {e}")
```

Listing 14 – Fonction log\_etl\_db() dans les scripts ETL

## A.3 — Flux de la sauvegarde

```
Azure SQL Database
|
|   Sauvegardes automatiques (complet + differentiel + logs)
|   gerees par Azure (PITR, retention 14 jours)
|
v
backup_to_datalake.py (etape 5 du pipeline ETL)
|
|   az sql db export --> fichier .bacpac
|
v
Data Lake ADLS Gen2
raw/backups/projet_data_eng_20260217_143000.bacpac
|
+-- Log dans dwh.log_etl
+-- Nettoyage automatique des .bacpac > 30 jours
```

Listing 15 – Flux de la sauvegarde

## A.4 — Export BACPAC vers le Data Lake (backup\_to\_datalake.py)

```
1 def export_bacpac(config, resource_group):
2     """Exporte la base Azure SQL en .bacpac via az sql db export."""
3     storage_account = config['storage_account']
4     timestamp = datetime.now().strftime('%Y%m%d_%H%M%S')
5     blob_name = f"backups/{config['database']}_{timestamp}.bacpac"
6
7     # Recuperer la cle du storage account
8     result = subprocess.run(
9         ['az', 'storage', 'account', 'keys', 'list',
```

```

10         '--account-name', storage_account,
11         '--resource-group', resource_group,
12         '--query', '[0].value', '-o', 'tsv'],
13         capture_output=True, text=True)
14     storage_key = result.stdout.strip()
15
16     storage_uri = (f"https://{storage_account}.blob.core.windows.net"
17                   f"/raw/{blob_name}")
18
19     # Exporter la base en .bacpac directement vers le blob
20     subprocess.run([
21         'az', 'sql', 'db', 'export',
22         '--admin-user', config['user'],
23         '--admin-password', config['password'],
24         '--name', config['database'],
25         '--server', config['server'],
26         '--resource-group', resource_group,
27         '--storage-key', storage_key,
28         '--storage-key-type', 'StorageAccessKey',
29         '--storage-uri', storage_uri,
30     ], capture_output=True, text=True, timeout=1800)
31
32     return blob_name

```

Listing 16 – Export BACPAC vers le Data Lake (backup\_to\_datalake.py)

## A.5 — Nettoyage automatique des anciens backups

```

1 def cleanup_old_backups(config, resource_group, retention_days=30):
2     """Supprime les .bacpac de plus de retention_days jours."""
3     result = subprocess.run(
4         ['az', 'storage', 'blob', 'list',
5          '--account-name', config['storage_account'],
6          '--container-name', 'raw', '--prefix', 'backups/',
7          '--query', '[0].{name:name}', '-o', 'tsv'],
8         capture_output=True, text=True)
9
10    for blob_name in result.stdout.strip().split('\n'):
11        # Extraire la date du nom (db_YYYYMMDD_HHMMSS.bacpac)
12        blob_date = datetime.strptime(...)
13        if (datetime.now() - blob_date).days > retention_days:
14            subprocess.run(['az', 'storage', 'blob', 'delete',
15                           '--account-name', config['storage_account'],
16                           '--container-name', 'raw',
17                           '--name', blob_name])

```

Listing 17 – Nettoyage automatique des anciens backups

## A.6 — Restauration via Azure CLI

```

1 # Methode 1 : Restauration Point-in-Time (PITR)
2 # Restaure la base a un instant precis (dans les 14 derniers jours)
3 az sql db restore \
4     --resource-group rg-projet-data \
5     --server sqllelbrek-prod \
6     --name projet_data_eng \

```

```

7      --dest-name projet_data_eng_restored \
8      --time "2026-02-16T10:00:00Z"
9
10     # Methode 2 : Import d'un fichier BACPAC depuis le Data Lake
11     # Recree la base a partir d'un export stocke dans le Data Lake
12 az sql db import \
13     --resource-group rg-projet-data \
14     --server sqllelbrek-prod \
15     --name projet_data_eng_restored \
16     --storage-key <storage_key> \
17     --storage-key-type StorageAccessKey \
18     --storage-uri "https://dlelbrek.blob.core.windows.net/raw/\
19 backups/projet_data_eng_20260217_143000.bacpac" \
20     --admin-user sqladmin --admin-password ***

```

Listing 18 – Restauration via Azure CLI

## A.7 — Tables du schema security (011\_security\_rls.sql)

```

1  -- 101 agences : villes > 10 000 hab dans la region
2  CREATE TABLE security.agences (
3      agence_id          INT          IDENTITY(1,1) PRIMARY KEY,
4      commune_code       NVARCHAR(10) NOT NULL,
5      ville              NVARCHAR(100) NOT NULL,
6      departement_code   NVARCHAR(3)  NOT NULL,
7      taille_agence      NVARCHAR(20) NOT NULL -- 'GRANDE', 'MOYENNE', '
      PETITE'
8  );
9
10 -- Employes avec hierarchie (519 au total)
11 CREATE TABLE security.employes (
12     employe_id          INT          IDENTITY(1,1) PRIMARY KEY,
13     login_sql           NVARCHAR(100) NOT NULL UNIQUE,
14     niveau_hierarchique NVARCHAR(30) NOT NULL,
15     agence_id           INT          NULL REFERENCES security.agences,
16     departement_code    NVARCHAR(3)  NULL, -- NULL = region entiere
17     manager_id          INT          NULL REFERENCES security.employes
18 );
19
20 -- Mapping login -> departements autorises (alimente par ETL)
21 CREATE TABLE security.utilisateurs_zones (
22     login_sql           NVARCHAR(100) NOT NULL,
23     departement_code    NVARCHAR(3)  NULL -- NULL = acces region entiere
24 );

```

Listing 19 – Tables du schema security (011\_security\_rls.sql)

## A.8 — Fonction predicat RLS (011\_security\_rls.sql)

```

1  CREATE FUNCTION security.fn_rls_geographie(@departement_code NVARCHAR
2      (3))
3  RETURNS TABLE WITH SCHEMABINDING AS
4  RETURN
5      SELECT 1 AS acces_autorise
6      WHERE
7          -- Utilisateur absent de utilisateurs_zones -> acces total

```

```

7      -- (role_admin, role_analyst, role_etl_process)
8      NOT EXISTS (
9          SELECT 1 FROM security.utilisateurs_zones
10         WHERE login_sql = USER_NAME()
11     )
12     OR
13     -- Utilisateur enregistre (role_consultant) -> filtrer par dept
14     EXISTS (
15         SELECT 1 FROM security.utilisateurs_zones
16         WHERE login_sql = USER_NAME()
17             AND (departement_code = @departement_code
18                 OR departement_code IS NULL) -- NULL = region
19                                     entiere
20     );
21
22 GO
23
24 CREATE SECURITY POLICY security.policy_rls_geographie
25 ADD FILTER PREDICATE
26     security.fn_rls_geographie(departement_code)
27 ON dwh.dim_geographie
28 WITH (STATE = ON, SCHEMABINDING = ON);

```

Listing 20 – Fonction predicat RLS (011\_security\_rls.sql)

## A.9 — Extrait de code 9

```

1 CREATE VIEW dm.vm_nouveau_datamart AS
2 SELECT ... FROM dwh.fait_xxx f
3 JOIN dwh.dim_xxx d ON f.xxx_id = d.xxx_id
4 GROUP BY ...;

```

## A.10 — Implementation SCD Type 1 – Ecrasement

```

1 CREATE PROCEDURE dwh.sp_scd_type1_activite
2     @code_naf NVARCHAR(10),
3     @nouveau_libelle NVARCHAR(200)
4 AS
5 BEGIN
6     -- Journaliser l'ancienne valeur (pour audit)
7     INSERT INTO dwh.log_etl (etape, table_cible, statut, message)
8     SELECT 'SCD_TYPE1', 'dim_activite', 'INFO',
9         'Ancien libelle: ' + libelle_section
10    FROM dwh.dim_activite
11    WHERE code_section_naf = @code_naf;
12
13    -- Ecrasement de la valeur
14    UPDATE dwh.dim_activite
15    SET libelle_section = @nouveau_libelle,
16        date_modification = GETDATE()
17    WHERE code_section_naf = @code_naf;
18 END;

```

Listing 21 – Implementation SCD Type 1 – Ecrasement

## A.11 — Procedure SCD Type 2 – Historisation

```
1 CREATE PROCEDURE dwh.sp_scd_type2_geographie
2     @code_commune    NVARCHAR(10),
3     @nouveau_nom     NVARCHAR(200),
4     @nouveau_code    NVARCHAR(10) = NULL
5 AS
6 BEGIN
7     DECLARE @ancienne_version INT
8     DECLARE @ancien_id INT
9
10    -- Recuperer la version courante
11    SELECT @ancienne_version = version, @ancien_id = geo_id
12    FROM dwh.dim_geographie
13    WHERE code_commune = @code_commune AND est_actif = 1;
14
15    -- Fermer l'ancien enregistrement
16    UPDATE dwh.dim_geographie
17    SET date_fin_validite = GETDATE(),
18        est_actif = 0,
19        date_modification = GETDATE()
20    WHERE geo_id = @ancien_id;
21
22    -- Insérer la nouvelle version
23    INSERT INTO dwh.dim_geographie (
24        code_departement, nom_departement, code_region,
25        nom_region, niveau, code_commune, nom_commune,
26        latitude, longitude,
27        date_debut_validite, date_fin_validite,
28        est_actif, version,
29        date_creation, date_modification
30    )
31    SELECT
32        code_departement, nom_departement, code_region,
33        nom_region, niveau,
34        ISNULL(@nouveau_code, code_commune),
35        @nouveau_nom,
36        latitude, longitude,
37        GETDATE(), NULL,
38        1, @ancienne_version + 1,
39        GETDATE(), GETDATE()
40    FROM dwh.dim_geographie
41    WHERE geo_id = @ancien_id;
42
43    -- Journaliser le changement
44    INSERT INTO dwh.log_etl (etape, table_cible, statut, message)
45    VALUES ('SCD_TYPE2', 'dim_geographie', 'SUCCES',
46        'Nouvelle version (v' + CAST(@ancienne_version + 1 AS VARCHAR)
47        + ') pour commune ' + @code_commune
48        + ' -> ' + @nouveau_nom);
49 END;
```

Listing 22 – Procedure SCD Type 2 – Historisation

## A.12 — Procedure SCD Type 3

```
1 CREATE PROCEDURE dwh.sp_scd_type3_demographie
```

```

2      @code_pcs          NVARCHAR(10),
3      @nouveau_libelle NVARCHAR(200)
4  AS
5  BEGIN
6      -- Sauvegarder l'ancien libelle et mettre a jour
7      UPDATE dwh.dim_demographie
8      SET ancien_libelle_pcs = libelle_pcs,
9          libelle_pcs = @nouveau_libelle,
10         date_changement_pcs = GETDATE(),
11         date_modification = GETDATE()
12      WHERE code_pcs = @code_pcs
13             AND libelle_pcs <> @nouveau_libelle;
14
15      INSERT INTO dwh.log_etl (etape, table_cible, statut, message)
16      VALUES ('SCD_TYPE3', 'dim_demographie', 'SUCCES',
17              'PCS ' + @code_pcs + ' mis a jour avec Type 3');
18  END;

```

Listing 23 – Procedure SCD Type 3

### A.3 Glossaire

Terme	Définition
CSV	Comma-Separated Values – Format de fichier texte pour les données tabulaires
CDC	Change Data Capture – Mécanisme de capture des changements dans une base de données
DWH	Data Warehouse – Entrepôt de données
ETL	Extract, Transform, Load – Processus d'extraction, transformation et chargement de données
ITIL	Information Technology Infrastructure Library – Référentiel de bonnes pratiques IT
KPI	Key Performance Indicator – Indicateur clé de performance
MCO	Maintien en Conditions Opérationnelles
MERGE	Instruction SQL combinant INSERT et UPDATE en une seule opération
RBAC	Role-Based Access Control – Contrôle d'accès basé sur les rôles
RLS	Row-Level Security – Mécanisme SQL Server filtrant les lignes selon l'identité de l'utilisateur, de manière transparente et non contournable
RGPD	Règlement Général sur la Protection des Données
RPO	Recovery Point Objective – Perte de données maximale acceptable
RTO	Recovery Time Objective – Temps de restauration maximal acceptable
SCD	Slowly Changing Dimension – Dimension à évolution lente

TABLE 7 – Glossaire des termes techniques

### A.4 Arborescence du projet

```

Projet-Data-ENG-E6/
  Terraform/
    main.tf                                # Infrastructure Azure + retention
    backups [E6]

```



```

adf_linked_services.tf      # Data Factory
rbac_configuration.tf       # RBAC Azure
variables.tf / outputs.tf
sql/
  001_create_schemas.sql   # 4 schemas (stg, dwh, dm, analytics)
  002_create_dimensions.sql # 6 dimensions
  003_create_facts.sql      # 7 tables de faits
  004_populate_dimensions.sql
  005_create_datamarts.sql  # 5 datamarts + dashboard
  006_configure_security.sql # [E6] 4 roles RBAC (refonte + suppr.
    anciens roles)
  007_configure_performance.sql
  008_configure_logging.sql  # [E6] log_etl, log_erreurs, vues
    monitoring
  009_configure_backup.sql   # [E6] Vues monitoring backups
  010_scd_dimensions.sql     # [E6] SCD Type 1/2/3, MERGE, vues SCD
  011_security_rls.sql       # [E6] Schema security, RLS, agences,
    employes
  deploy_dwh.py
analytics/
  lib/
    data_prep.py            # [E6] Fonctions partagees ETL (
      load_communes, etc.)
  etl/
    run_etl.py              # Orchestrateur ETL [E6: logging +
      backup]
    load_dimensions.py       # Dimensions + SCD Type 1/2/3 [E6]
    load_facts.py            # Faits + emploi + menages [E6]
    backup_to_datalake.py    # [E6] Export BACPAC vers ADLS Gen2
    load_security.py         # [E6] Chargement agences, employes,
      zones RLS
    etl_notifier.py          # [E6] Notifications email (succes +
      erreur)
  tests/
    test_e6_evolution.py     # 10 tests E6 (journalisation, backup,
      sources, SCD, RBAC)
docs/
  E5_DATA_WAREHOUSE_GUIDE.md
  E6_MAINTENANCE_METHODODOLOGY.md # [E6] ITIL, incidents, KPIs
  E6_SCALABILITY_PROCEDURES.md   # [E6] Procedures d'evolution
rapports/
  E6_Evolution_Entrepot/
    rapport_E6.tex           # Ce rapport
uploads/
  landing/csv/
    EMPLOI_CHOMAGE_hauts_de_france.csv # Source emploi [E6]
    Menage_hauts_de_france.csv         # Source menages [E6]

```