

Moderniser un projet *dbt* legacy

Data  
Days

Tests, observabilité  
et migration sans risque



=elementary

[sf≡ir]



Henri-Maxime  
**Ducoulombier**

Senior Data Engineer

[ducoulombier.hm@sfeir.com](mailto:ducoulombier.hm@sfeir.com)



SÉCURITÉ — GPG — BONNES PRATIQUES

## Sécuriser l'échange de données sensibles avec GPG

Chaque fois qu'un mot de passe ou un fichier sensible est partagé par mail, un chaton meurt quelque part dans le monde. Ça n'est pas une fatalité et il existe une solution pour stopper cette tragédie: le chiffrement des échanges avec GPG (Gnu Privacy Guard)



DATA — DBT — BONNES PRATIQUES — CODE — TIPS

## DBT : des macros en cascade

La gestion des packages dbt avec l'override des macros génériques et par domaine.



SFEIR! DEV



AUTOMATION TEST — TIPS — BONNES PRATIQUES — SUCCESS STORY

## Vieille flûte!

L'histoire vraie et incroyable de l'affaire "Vieille flûte!".

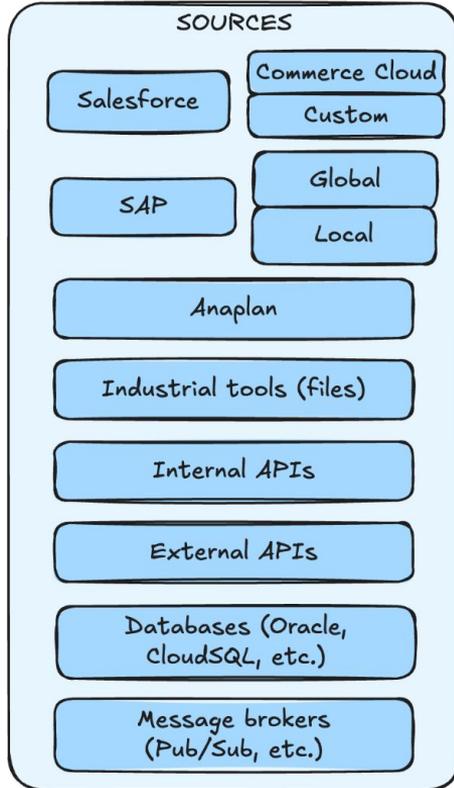
HENRI-MAXIME DUCCOULOMBIER

21 février 2024 · 11:53 AM — 6 lecture min

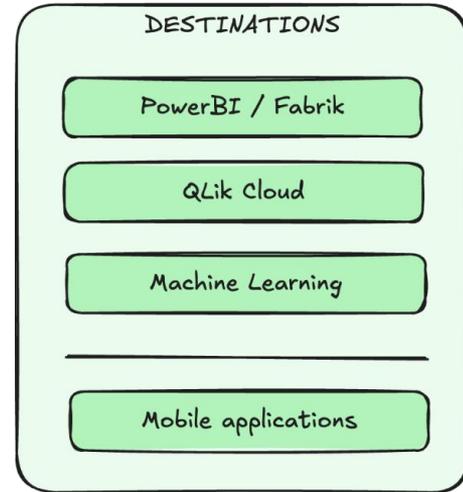
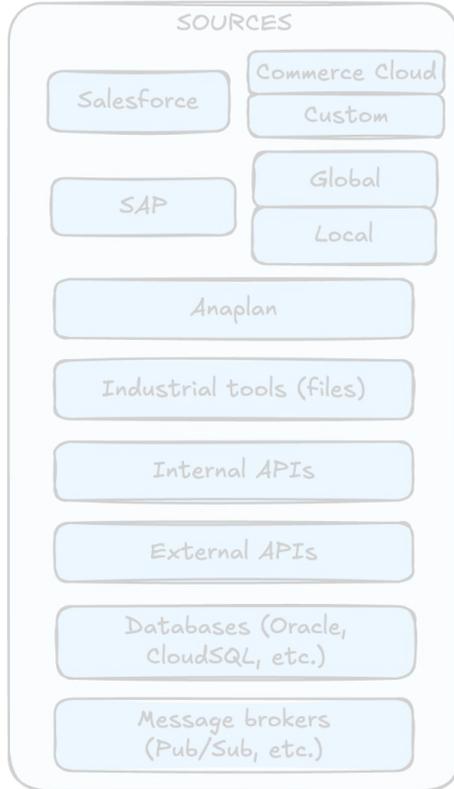


**L'architecture et la stack data actuelle**

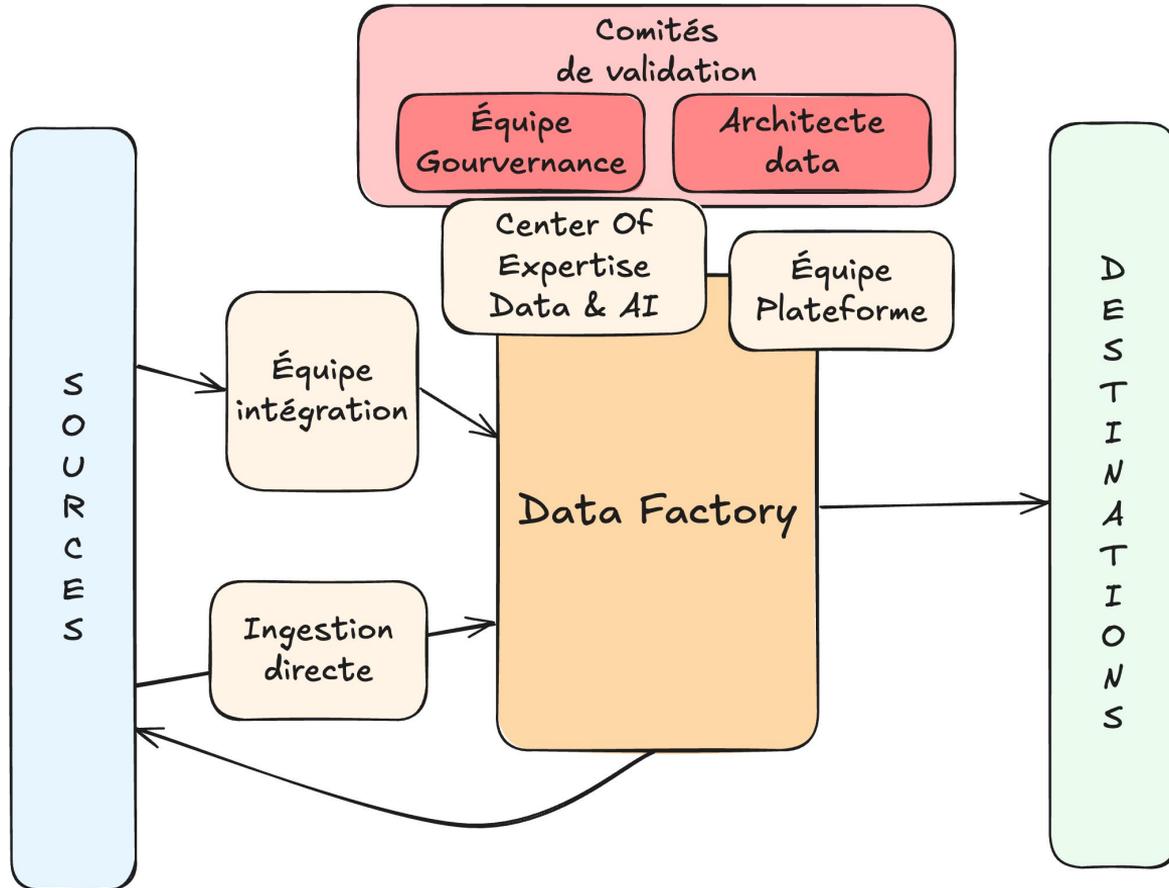
# Architecture globale → sources



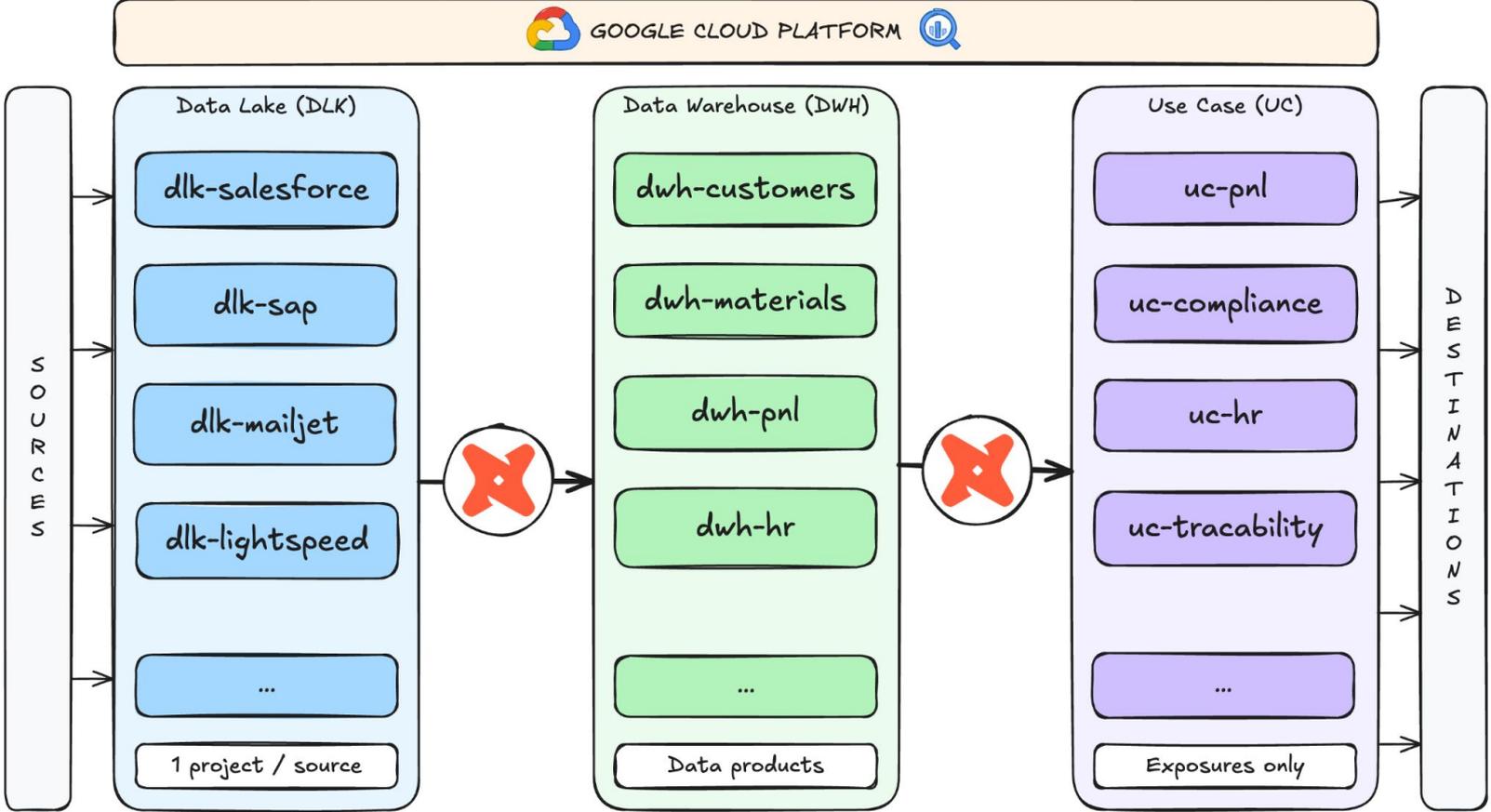
# Architecture globale → destinations



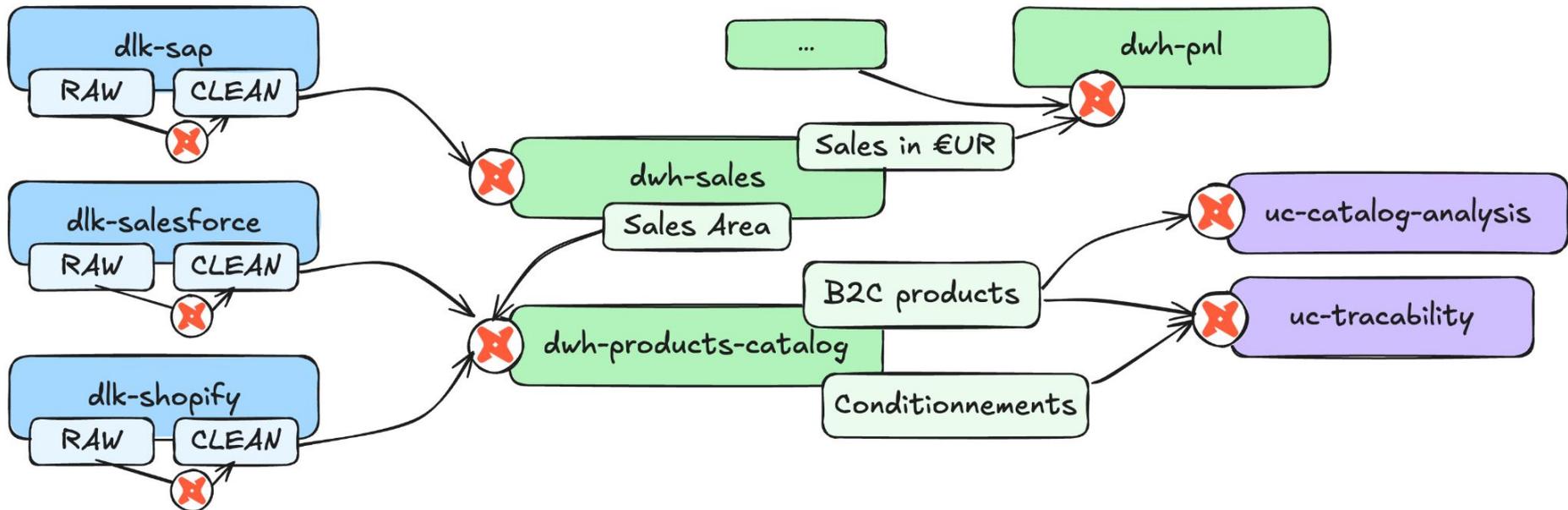
# Architecture globale → data factory



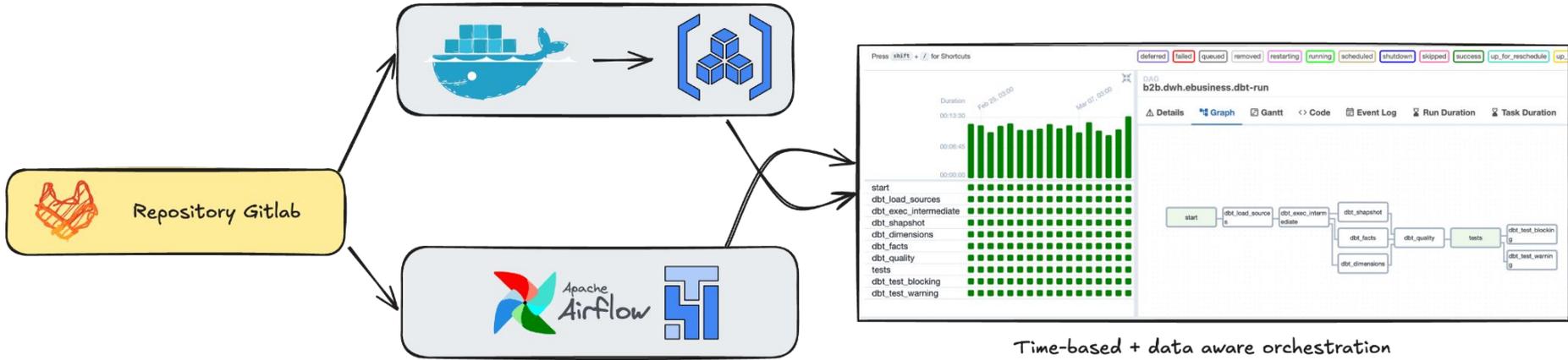
# DLK → DWH → UC



# Dépendances entre les 3 layers

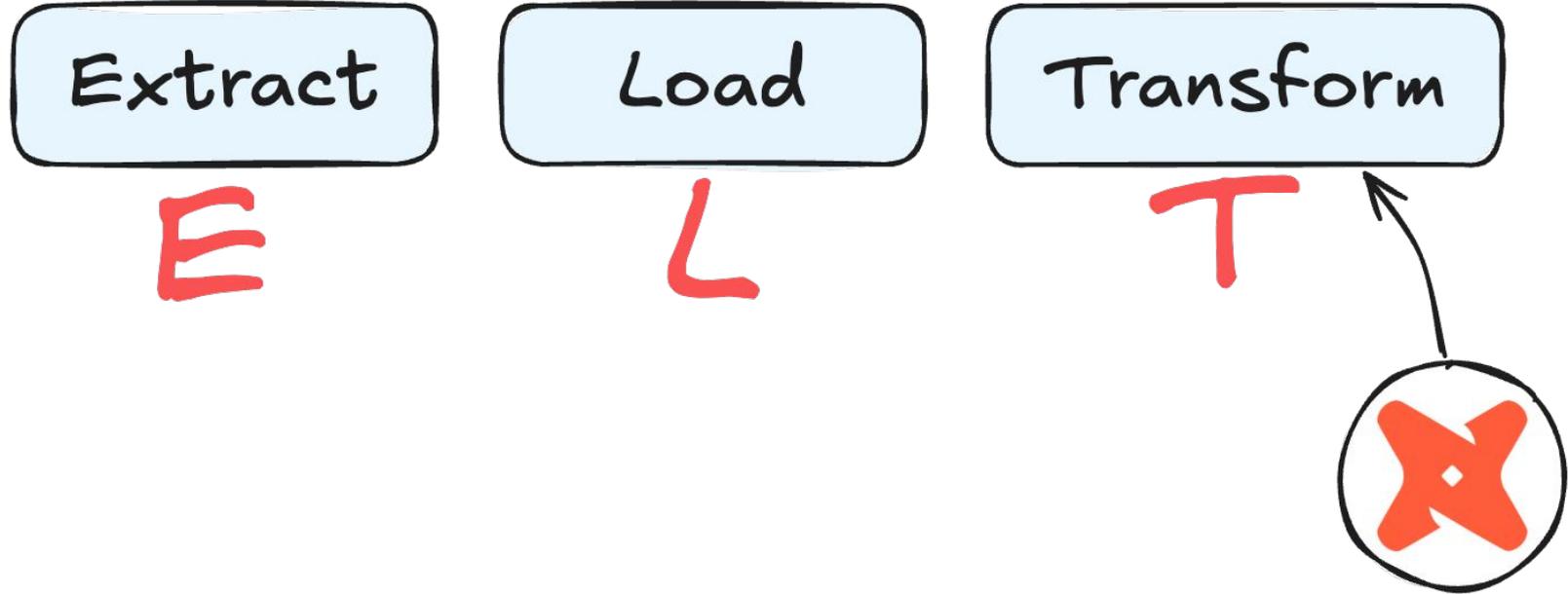


# Orchestration des tâches

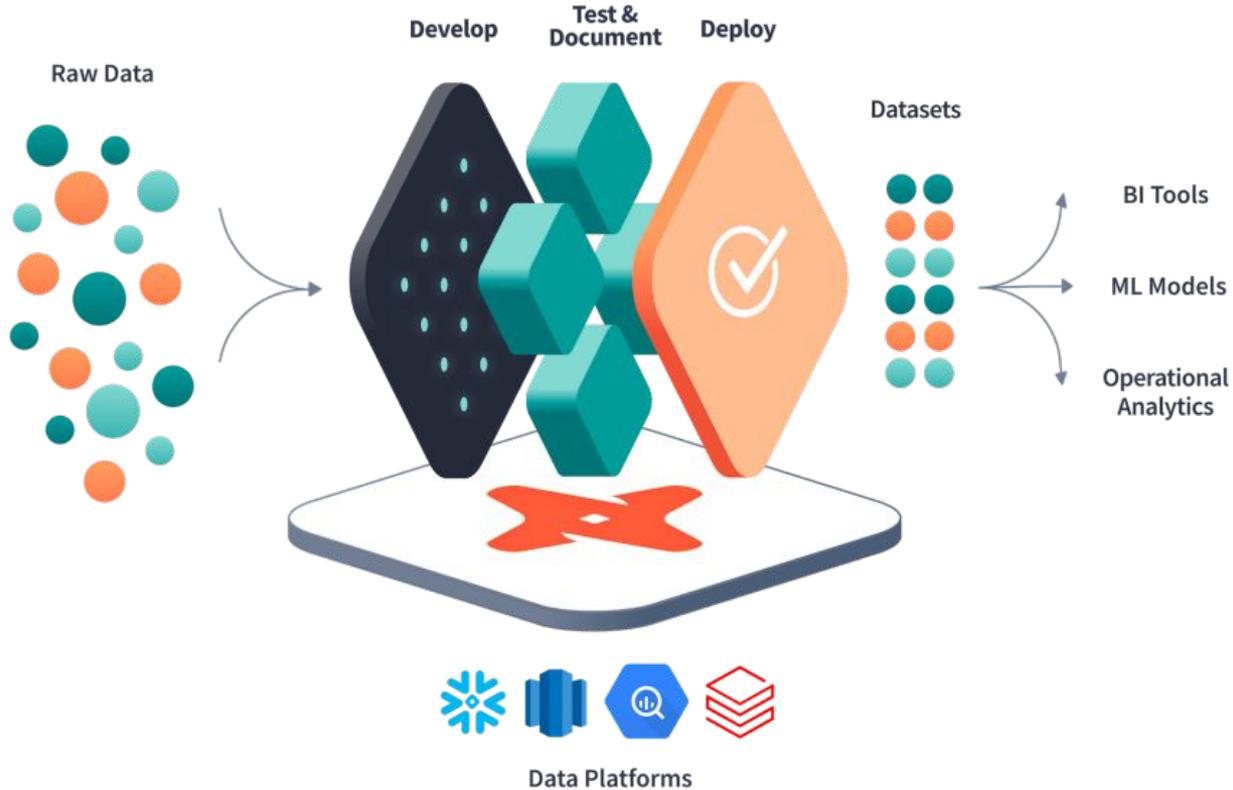




# Data Build Tool



# Principe de fonctionnement



## **dbt c'est...**

- **Transformations de données décrites en SQL**
- **Références et dépendances automatiques entre modèles**
- **Versionning du code et collaboration facilitée**
- **Tests, documentation et lineage**
- **Contrats et versions de modèles**
- **Modularité et personnalisation**

# Killer feature: tests

## Data tests

- Valider la donnée dans les sources
- Valider la donnée une fois transformée
- Permet d'interrompre le pipeline en cas de problème

## Unit tests

- Valider le code en dev et dans la CI/CD
- Valider des règles de calcul, du code métier
- Permet de ne pas déployer du code défectueux

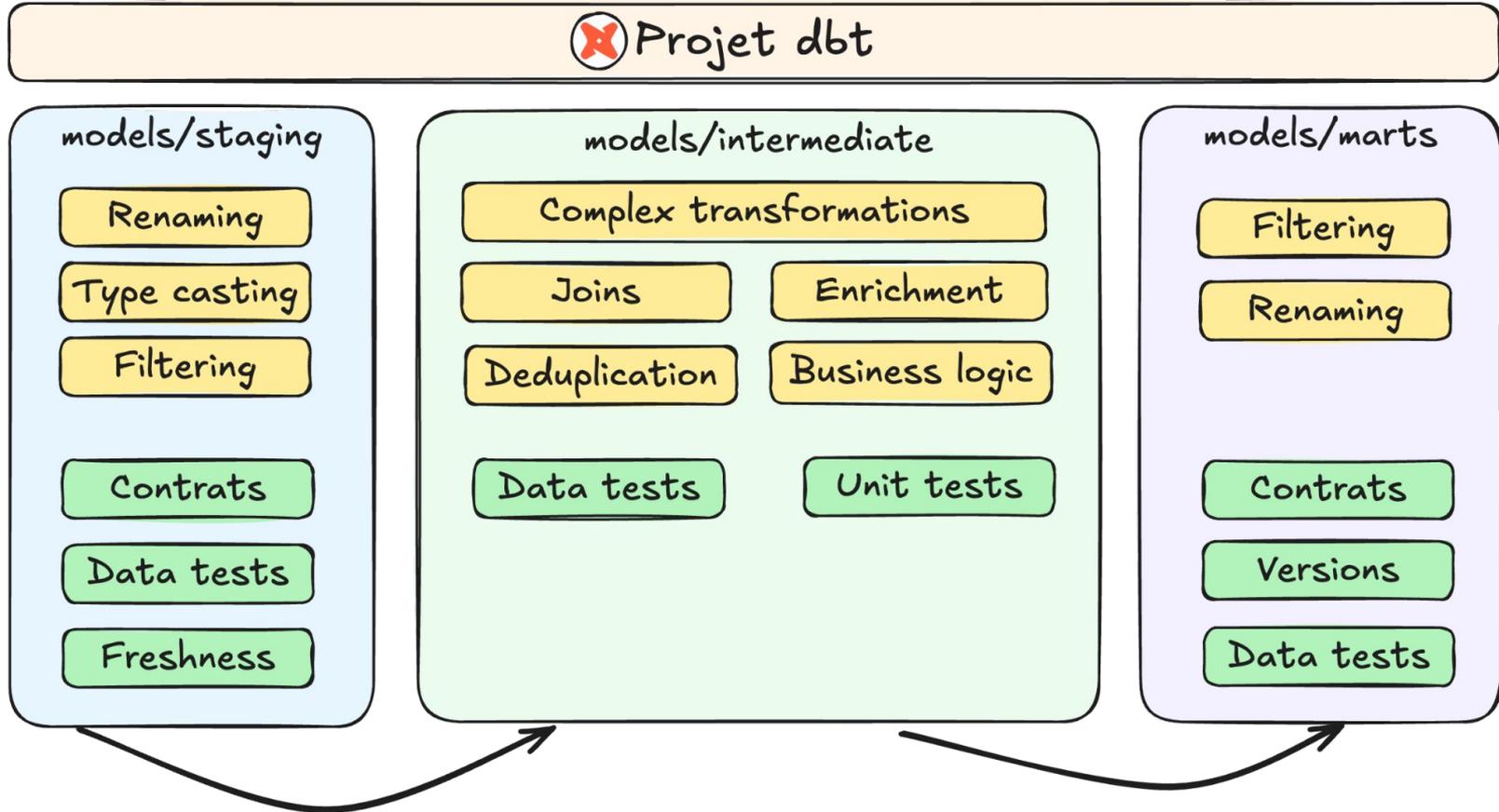
## **Data contracts**

- **Garantir la validité** des schémas de données
- **Ajouter des contraintes, déclarer des clés**

## **Versions de modèles**

- **Gérer les *breaking changes*** dans les modèles
- **Introduire la notion de cycle de vie** des modèles

# Organisation d'un projet dbt



**Rome ne s'est pas faite en un jour...**

# De vieux projets...

## Premiers projets *dbt* datent de 2021

- Version de dbt actuelle : **dbt 1.9** (1.10 en beta)
- Early adopters = essayer les plâtres
- Ancien projet "*dataplatfom*" qui ne devrait plus exister

```
docker pull fishtownanalytics/dbt:0.21.1
```

# Arrivée tardive de la gouvernance

Mise en place réelle depuis début 2024

- *Policies difficiles à appliquer sur l'existant*
- *Beaucoup + de contraintes qu'auparavant*

# Empilement et duplication de données

## Difficulté à cataloguer et cartographier

- *Beaucoup de données, beaucoup de doublons*
- *Les politiques historiques pas toujours respectées*
- *Shadow IT*

## Criticité des applications existantes

- *Nécessité absolue de maintenir les use cases utilisés*
- *Pas de budget pour tout refaire*
- *Pas de volonté de tout refaire*

**Error in SQL line 3871**



# La règle des 3U



UTILE



UTILISÉ



UTILISABLE

## Principaux problèmes sur le périmètre:

- *Non respect de l'architecture DLK -> DWH -> UC*
- *Vieilles versions des outils et packages*
- *Non respect des patterns et templates CI/CD*
- *Code mort caché un peu partout*
- *Opérations manuelles régulières (RUN x 1000)*

# Ne pas nourrir le monstre

## Besoins d'évolution

- *Supporter les nouveautés des sources existantes*
- *Ajouter de nouvelles sources de données*
- *Améliorer / corriger les KPIs existants*
- *Ajouter de nouveaux KPIs*
- *Appliquer de la compliance réglementaire*

## Comment procéder ?

- **Étape n°1 : Identifier le code mort et le désactiver**

- Ne pas migrer ce qui est inutile
- Transformer ce qui est inutilement complexe

- **Étape n°2 : Migration technique des outils**

- Uniformiser les versions de *dbt*
- Corriger les configurations (auth, packages, etc.)
- Mettre en place *elementary data sur le legacy*

- **Étape n°3 : Réorganisation de la couche DLK**
  - Appliquer le pattern "1 source = 1 dlk"
  - Mise en place de tests et de contrats *dbt*
- **Étape n°4 : Réorganisation des projets *dbt DWH et UC***
  - Respecter la structure recommandée *stg / int / mart*
  - Mise en place de tests data / unitaires
- **Étape n°5 : Damage control**

**Et elementary dans tout ça ?**

# Constat et besoins

- **Les journaux ne sont pas suffisants**
- **Besoin technique et métier**
- **Intégration forte avec *dbt***
- **Compatibilité avec tous types de tests**
- **Si possible avec une interface sympa**

# Le candidat idéal



+

**=elementary**

# Fonctionnement d'*elementary*

- Package *python* ET package *dbt*
- *elementary* surcharge TOUS les tests *dbt*
- **pre-hook** et **post-hook** à chaque exécution *dbt*
- Génération asynchrone des rapports
- Possibilité d'envoyer des alertes

# Setup d'elementary

- À l'installation et à chaque montée de version  
→ Init / upgrade des modèles *elementary*

```
dbt run --select elementary
```

- **Création obligatoire d'un profile *dbt* dédié**  
→ L'outil cherche un profil nommé "elementary"
- **Ajout d'un "selector" par défaut qui ignore *elementary***
  - Se fait dans le *selectors.yml*
    - Ne pas lancer le setup à chaque fois
  - Désactiver les hooks facilement via une variable d'env

# Configuration de l'orchestration

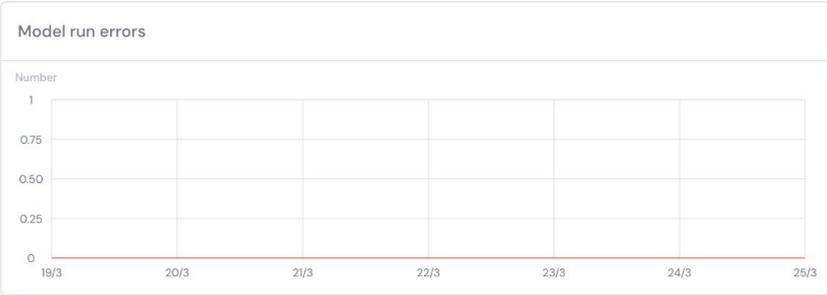
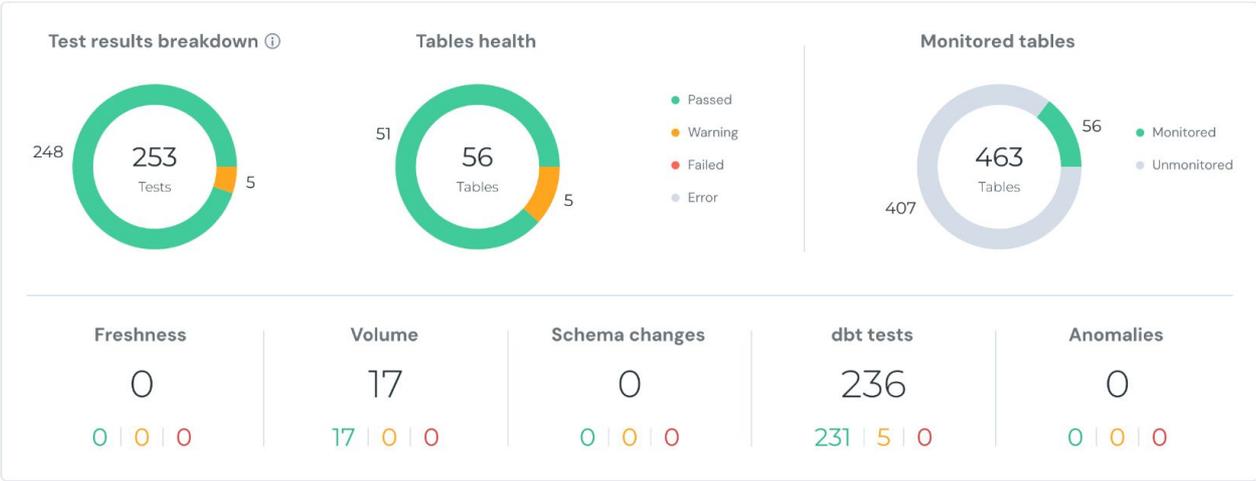
- *elementary* qui plante → non critique
- **générer les rapports = longue opération**  
→ le faire en asynchrone
- **génération de rapport peu importe le statut**  
→ rapport d'erreurs / warnings

# Exemple de dashboard – homepage

- Dashboard**  
Review data health
- Lineage**  
Analyze root cause & impact
- Catalog**  
Explore datasets
- DATA TESTS**
- Results**  
Triage data issues
- Test Execution History**  
Find frequently failing tests
- Test Configuration**  
Add, delete & edit tests
- DATA PIPELINES**
- Model Duration**  
Detect pipeline bottlenecks
- INTEGRATIONS**
- Environments**  
Connect databases, alerts & BI

## Dashboard

Filter by [dropdown] [dropdown] Last 7 Days



# Exemple de dashboard – tests results

Data Tests > Results

Last 7 Days

**All Tables** 248 5

Owners: -    Tags: -    Path: dbt\_project/

Table name	Column name	Test name	Test type	Last test run	Last status
product_target_usage	product_guid	relationships	generic	2025-03-25 04:01:18	Warning

relationships [View test runs](#) [View in lineage](#) [Copy link](#)

from_field
1588AF8A-E422-4D50-9F8B-7ED04B6D21B4
2123700B-1145-478B-883C-204B838F7A64
6385756E-75E2-48C7-8747-A82849651B5C
83A32FA0-F5DD-41A1-9B6B-B8DBE3F7475B
83A32FA0-F5DD-41A1-9B6B-B8DBE3F7475B

**Description**

Test Description:  
This test validates that all of the records in a child table have a corresponding record in a parent table. This property is referred to as "referential integrity".

**Result**

Result Description:  
Got 21 results, configured to warn if != 0

Result Query:  
with child as ( select product\_guid as from\_fi...

**Configuration**

Test Name:  
relationships

Test Params:  
{ "to": "ref('PRODUCT')", "field": "product\_guid", "colu...

# Exemple de dashboard – temps d'exécution

Data Pipelines > Model Duration

Last 7 Days

## All Tables

Owners: -    Tags: -    Path: dbt\_project/

Schema	Name	Last exec. time	Median exec. time	Ex time change rate	Results	Last run status	
dprd_facturation	contrat_ligne_prestatio...	10.3 Secs	9.1 Secs	<span style="color: red;">↗ 13.2%</span>	126	Success	-
Model execution time and status over time							<a href="#">View in lineage</a> <a href="#">Copy link</a>
dprd_facturation	contrat_ligne	8.9 Secs	9.2 Secs	<span style="color: green;">↘ 3.3%</span>	126	Success	+
dprd_facturation	contrat	7.7 Secs	7.9 Secs	<span style="color: green;">↘ 2.5%</span>	126	Success	+
dprd_facturation	realise	7 Secs	6.7 Secs	<span style="color: red;">↗ 4.5%</span>	126	Success	+

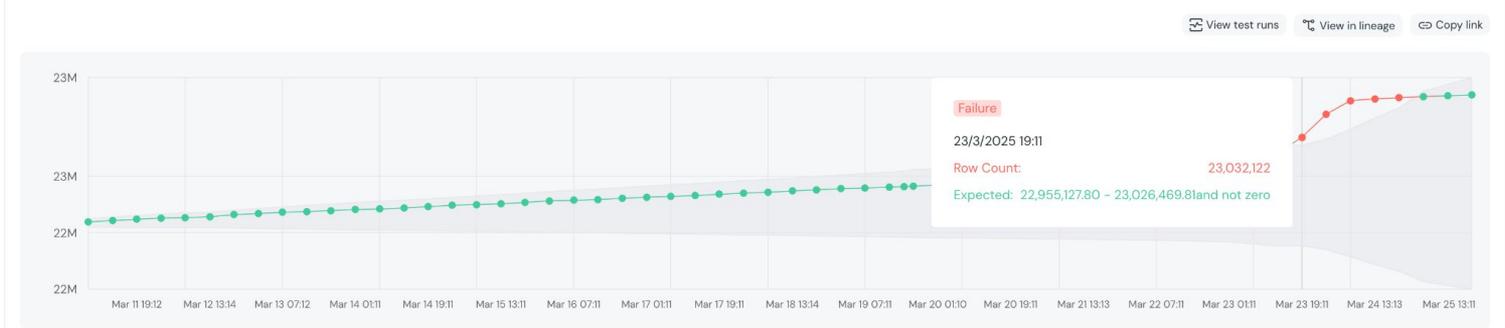
# Exemple de dashboard – anomalie de volume

Data Tests > Results

Last 7 Days

All Tables 28 2  
Owners: - Tags: - Path: dbt\_project/

Table name	Column name	Test name	Test type	Last test run	Last status
contacts_raw_v1	-	volume_anomalies	row_count	2025-03-25 13:13:00	Warning



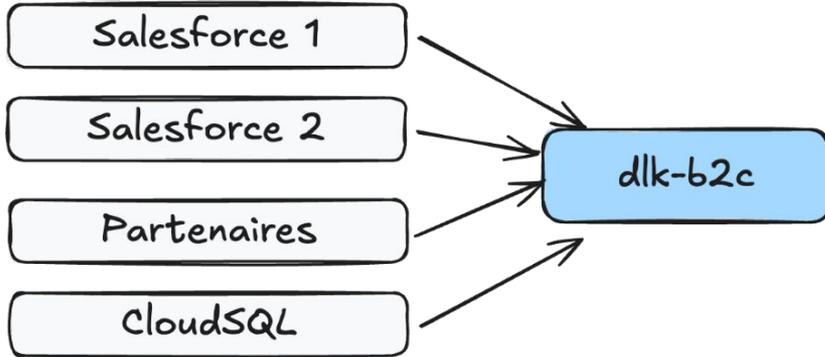
<b>Description</b> Test Description: Monitors the row count of your table over time.	<b>Result</b> Result Description: The last row_count value is 23062101. The average for this metric is 22999406.217.  Result Query: select * from (None) results where anomaly_score is not null and upp...	<b>Configuration</b> Test Name: volume_anomalies Timeframe: 1 day  Timestamp Column: No timestamp column Anomaly Threshold: 3
--------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------

raw_contacts	-	volume_anomalies	row_count	2025-03-25 13:04:20	Warning
--------------	---	------------------	-----------	---------------------	---------

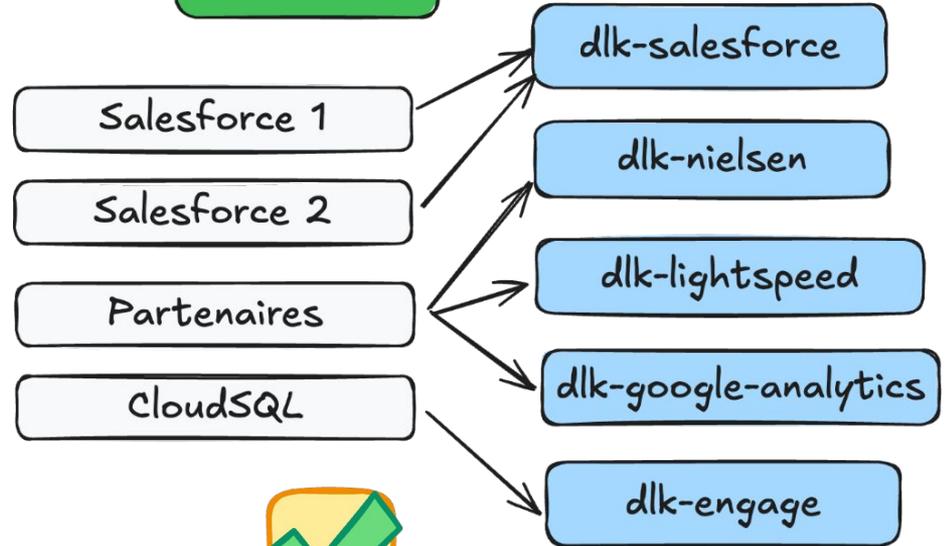
**Les changements appliqués**

# Réorganiser les sources

AVANT



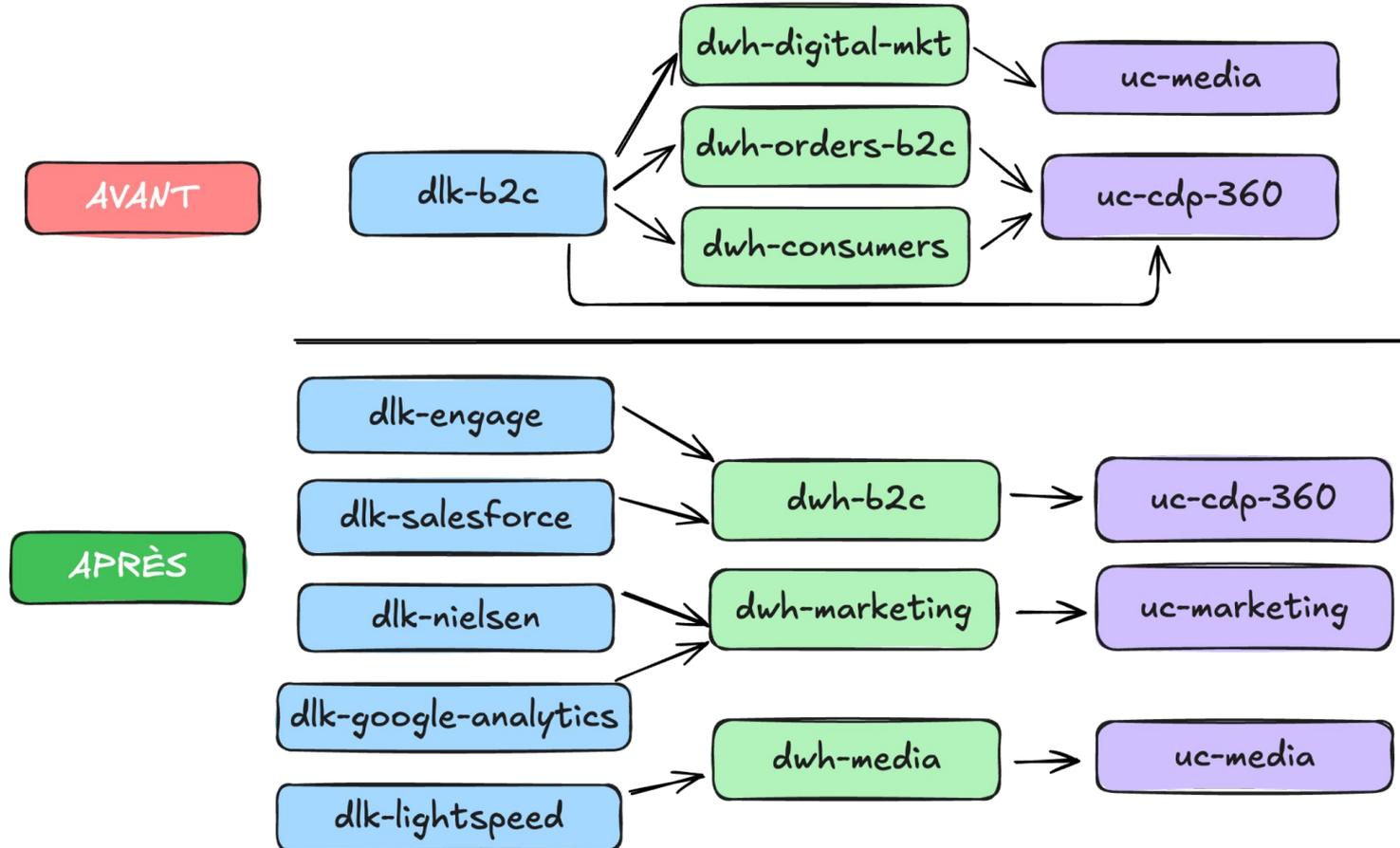
APRÈS



# Et tester les sources

- Tests sur volumes / dimensions / colonnes *=elementary*
- Mise en place de contrats et contraintes
- Si possible, ajout de timestamps et freshness

# Créer de nouveaux data products



# Et tester ces nouveaux data products

- Tests de réconciliation avant / après
- Data & Units Tests sur le code des règles métiers
- Contrats sur tous les datasets de sorties
- Mise en place de versions
- Suivi des temps d'exécution des pipelines

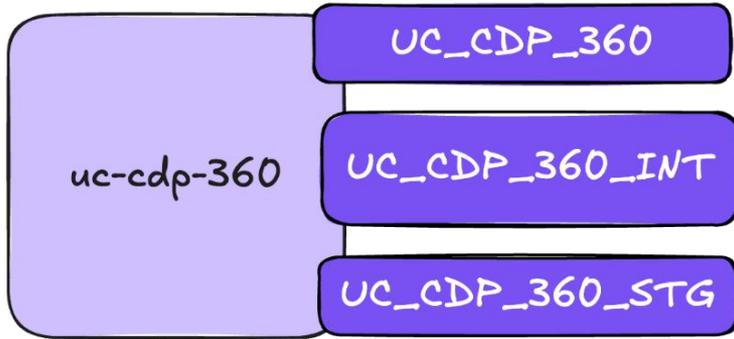
*=elementary*

*=elementary*

*=elementary*

# Exposer de nouveaux assets

AVANT



APRÈS



# Et tester ces nouveaux assets

- Mesure des volumes et distribution de donnée
- Contrats sur tous les datasets produits
- Mise en place de versions

*≡ elementary*

**Exemples de test**

# Les premiers tests basiques

## Data tests *génériques*

- Sur les colonnes
- Dans le YAML
- Via des packages
  - *dbt\_utils*
  - *dbt\_expectations*
  - *etc.*

```
1 models:
2   - name: INT__CONSUMERS
3     description: |
4       Merged sources for whole contact sources
5     columns:
6       - name: __source
7         description: |
8           Source system
9         data_tests:
10          - accepted_values:
11              values: ["SFCC", "SAPECC", "SAPB1", "SHOPIFY"]
12       - name: customer_code
13         description: |
14           Normalized customer code
15         data_tests:
16          - dbt_expectations.expect_column_values_to_match_regex:
17              regex: "^\\d{10}$"
```

# Les premiers tests basiques

## Data tests *génériques*

- Sur les sources & modèles
- Dans le YAML
- Via des packages
  - *dbt\_utils*
  - *dbt\_expectations*
  - *etc.*

```
1 sources:
2   - name: SAP
3     project: "dlk-sap-{{ env_var('ENV', 'dev') }}"
4     dataset: SAP"
5     tables:
6       - name: AUFK
7       - name: MAKT
8         data_tests:
9           - dbt_utils.unique_combination_of_columns:
10             combination_of_columns:
11               - MANDT
12               - MATNR
13               - SPRAS
14       - name: T001L
15       - name: T001W
```

# Gestion de la sévérité

## Error ou Warn ?

- **Error** → arrête le pipeline
- **Warn** → logs uniquement
- **Avec conditions**

```
1 models:
2   - name: INT__CONSUMERS
3     columns:
4       - name: __source
5         data_tests:
6           - accepted_values:
7             values: ["SFCC", "SAPECC", "SAPB1", "SHOPIFY"]
8             config:
9               severity: error
10              error_if: ">10000"
11              warn_if: ">100"
12   - name: phone_number
13     data_tests:
14       - not_empty:
15         config:
16           severity: warn
```

# Les génériques personnalisés

## Tests personnalisés

- **Sur les sources & modèles**
- **En SQL + YAML**
- **À écrire soit même**
- **Plusieurs tests en un seul**
- **Créer des packages**

```
1 {% test acme_phone_number_format(model, column_name, min_length=10) %}  
2 SELECT {{column_name }}  
3 FROM {{ model }}  
4 WHERE  
5     {{ column_name }} IS NULL  
6     OR LENGTH({{ column_name }}) < {{ min_length }}  
7     OR NOT REGEXP_CONTAINS({{ column_name }}, '^\\+?[0-9]+$')  
8 {% endtest %}
```

```
1 models:  
2   - name: INT__CONSUMERS  
3     columns:  
4       - name: mobile_phone_number  
5         data_tests:  
6           - acme_phone_number_format:  
7             minimum_length: 8
```

# Les tests singular

## Tests singular + complexes

- **Sur les sources & modèles**
- **En SQL uniquement**
- **Tests fonctionnels**
- **Utiliser les relations**

```
1 {{ config(severity = 'warn') }}
2
3 SELECT
4     o.order_id
5 FROM
6     {{ ref('orders') }} o
7 LEFT JOIN
8     {{ ref('order_items') }} oi ON o.order_id = oi.order_id
9 LEFT JOIN
10    {{ ref('products') }} p ON oi.product_id = p.product_id
11 GROUP BY
12     o.order_id
13 HAVING
14     COUNT(DISTINCT oi.product_id) > 10
15     OR SUM(oi.quantity * p.price) < 0
```

# Les tests de fraîcheur

## Avec ou sans timestamp

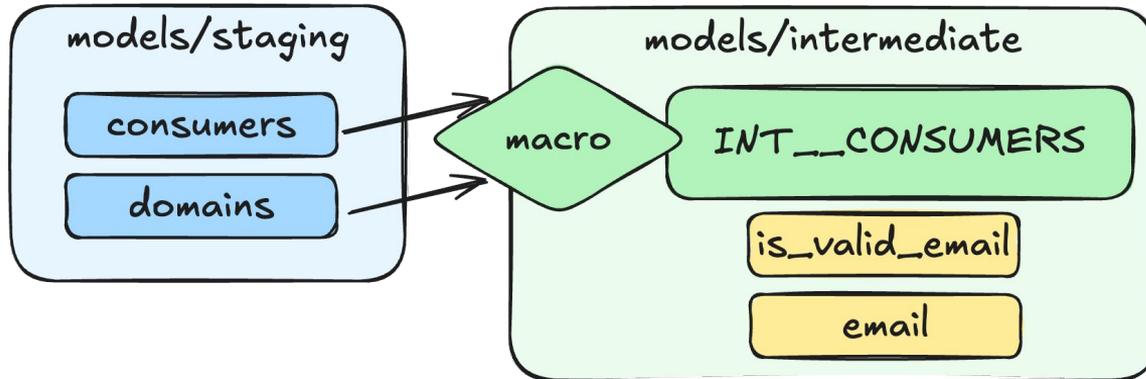
- Seulement sur les sources
- **Error** ou **Warn**
- En début de pipeline

```
sources:  
- name: SAP  
  project: "dlk-sap-{{ env_var('ENV', 'dev') }}"  
  dataset: SAP  
  loaded_at_field: __ingestionTimestamp  
  tables:  
    - name: MAKT  
      freshness:  
        warn_after: { count: 24, period: hour }  
        error_after: { count: 96, period: hour }  
    - name: T001W  
      loaded_at_field: null  
      freshness:  
        error_after: { count: 12, period: hour }
```

# Les tests unitaires

## Éviter les régressions

- Exécution locale et/ou dans la CI/CD
- Valider le code métier
- Long à écrire



# Les tests unitaires

## SI (given)

consumers

```
= [
  {"id": "101", "email": "xxx@yyy.com"},
  {"id": "102", "email": null},
  {"id": "103", "email": "zzz@xxx.com"}
]
```

domains

```
= [
  {"domain": "yyy.com"},
  {"domain": "yahoo.com"}
]
```

## ALORS (expect)

INT\_\_CONSUMERS

```
= [
  {"id": "101", "email": "xxx@yyy.com", "is_valid_email": true},
  {"id": "102", "email": null},
  {"id": "103", "email": "zzz@yyy.com"}
]
```

## SINON



**Contrats et versions**

# Contrats de données

## Forcer types et contraintes

- Au niveau *staging*
- Au niveau *marts (output)*
- **YAML** ← → **SQL**
- **avant l'écriture des données**

```
1 models:
2   - name: STG__EDA_REFERENTIALS
3     config:
4       contract:
5         enforced: true
6       constraints:
7         - type: primary_key
8           columns: [ reference_code, __filename ]
9     columns:
10      - name: __source
11        data_type: STRING
12      - name: reference_code
13        data_type: INT64
14      - name: referential_type_label
15        data_type: STRING
16      - name: reference_label
17        data_type: STRING
18      - name: __filename
19        data_type: STRING
20      - name: __ingestionTimestamp
21        data_type: FLOAT64
```

# Versions de modèles

## Avancer sans casser

- Au niveau *marts* (output)
- 1 fichier .sql par version
- Bien gérer la durée de vie
- Attention au coût

```
1 models:
2   - name: CONSUMERS
3     config:
4       contract:
5         enforced: true
6       latest_version: 2
7     versions:
8       - v: 1 # CONSUMERS_v1.sql
9         deprecation_date: "2024-10-17 00:00:00.00+00:00"
10        config:
11          enabled: false
12       - v: 2 # CONSUMERS_v2.sql
13         columns:
14           - name: preferences
15             data_type: ARRAY
16           - name: preferences.key
17             data_type: STRING
18           - name: preferences.value
19             data_type: ARRAY<STRING>
```

**Les tests *elementary***

# Elementary : paramétrage

## Quelques paramètres communs:

- *timestamp\_column* : analyser les changements dans le temps
- *severity* : réglable par test, **WARN** ou **ERROR**
- *anomaly\_direction* : détection des pics, des creux, ou des deux
- *detection\_period* et *training\_period*

# Elementary : tests de volume

```
1 models:
2   - name: CONTACTS_PROCESSED
3     config:
4       contract:
5         enforced: true
6       elementary:
7         timestamp_columns: updated_at
8     data_tests:
9       - elementary.volume_anomalies:
10         severity: warn
11         fail_on_zero: true
12       - at_least_one_not_null:
13         combination_of_columns:
14           - email
15           - wechat_id
16           - line_id
```

# Elementary : tests de répartition des valeurs

```
1 models:
2   - name: CONTACTS_PROCESSED
3     config:
4       contract:
5         enforced: true
6       elementary:
7         timestamp_columns: updated_at
8     data_tests:
9       - elementary.dimension_anomalies:
10         dimensions:
11           - __source
12           - brand
13           - market
```

# Elementary : anomalies de colonnes

```
1 models:
2   - name: CONTACTS_PROCESSED
3     config:
4       contract:
5         enforced: true
6       elementary:
7         timestamp_columns: updated_at
8     columns:
9       - name: first_order_delay
10      data_tests:
11        - elementary.column_anomalies:
12          severity: warn
13          column_anomalies:
14            - average
15            - standard_deviation
16            - min
17            - max
```

# Elementary : (toutes) anomalies de colonnes

```
1 models:
2   - name: CONTACTS_PROCESSED
3     config:
4       contract:
5         enforced: true
6       elementary:
7         timestamp_columns: updated_at
8     data_tests:
9       - elementary.all_columns_anomalies:
10         severity: warn
11         column_anomalies:
12           - average
13           - standard_deviation
14           - min
15           - max
```

**TL;DR**

# Leçons apprises



Respecter les guidelines,  
les polices



Utiliser les tests, les  
contrats, les versions



Travailler avec le métier  
pour avancer par étape



Restructurer au  
fur et à mesure



Décomposer les modèles  
en étapes



On oublie vite la syntaxe  
des insert / update



**MERCI !**