

# DW #2 — Fabrication (OEE) : de RAW → staging → dimensions → faits → OEE + Extensions

Programme Data Engineer

ECE

23 octobre 2025

# Objectifs du Mini-TP et plan

## Objectifs :

- Simuler des données **RAW** de fabrication : machines, produits, shifts, production, arrêts (downtime).
- Nettoyer/normaliser (*staging*) et construire les **dimensions** (date, machine, produit, shift).
- Construire deux tables de **faits** : `fact_production` (additive) et `fact_downtime` (semi-additive).
- Calculer l'**OEE** (Availability, Performance, Quality) par mois et tracer le graphique.

**Plan** : Contexte → RAW → Staging → DIM → FACT → Agrégations/OEE → Graphique → TP  
→ **Extensions (SCD2, Accumulating Snapshot, Bridge, LAF, DQ)**.

# Contexte et setup (imports, reproductibilité)

```
1 import numpy as np, pandas as pd
2 from datetime import datetime, timedelta
3 from pathlib import Path
4 import matplotlib.pyplot as plt
5
6 rng = np.random.default_rng(7)
7 out_dir = Path("dw_mfg"); out_dir.mkdir(exist_ok=True)
8
9 start, end = datetime(2025,1,1), datetime(2025,9,30)
10 dates = pd.date_range(start, end, freq="D")
```

# Génération RAW : référentiels (machines, produits, shifts)

```
1 areas = ["Cutting","Assembly","Packaging"]
2 raw_machines = pd.DataFrame({
3     "machine_id": range(1, 16),
4     "machine_name": [f"M{i:02d}" for i in range(1,16)],
5     "area": rng.choice(areas, 15),
6     "install_date": pd.to_datetime("2022-01-01") + pd.to_timedelta(rng
7         .integers(0,365,15), unit="D")
8 })
9 raw_products = pd.DataFrame({
10     "sku": [f"SKU{i:03d}" for i in range(1,41)],
11     "product_name": [f"Prod {i:03d}" for i in range(1,41)],
12     "ideal_rate_per_min": rng.choice([0.8,1.0,1.2,1.5], 40)    # u n i t s
13     /min
14 })
15 raw_shifts = pd.DataFrame({
16     "shift_id": [1,2,3],
```

# Génération RAW : production et arrêts (downtime)

```
1 # Production journalière par machine + shift + produit
2 prod_rows = []
3 for d in dates:
4     for m in raw_machines["machine_id"]:
5         for s in raw_shifts["shift_id"]:
6             if rng.random() < 0.85: # shift réellement produit ?
7                 sku = rng.choice(raw_products["sku"])
8                 good = rng.integers(200, 800)
9                 scrap = rng.integers(0, int(good*0.15))
10                prod_rows.append([d, m, s, sku, good, scrap])
11 raw_production = pd.DataFrame(prod_rows, columns=[
12     "prod_date", "machine_id", "shift_id", "sku", "units_good", "
13     units_scrap"
14 ])
15 # Arrêts (downtime) journaliers
16 failure_modes = ["Setup", "Breakdown", "Quality", "Material", "
    PlannedMaintenance"]
```

## Staging : nettoyage standardisation

```
1 stg_machines = raw_machines.assign(  
2     machine_name = raw_machines["machine_name"].str.upper(),  
3     area = raw_machines["area"].str.title()  
4 )  
5  
6 stg_products = raw_products.assign(  
7     product_name = raw_products["product_name"].str.title()  
8 )  
9  
10 stg_shifts = raw_shifts.copy()  
11  
12 stg_production = raw_production.assign(  
13     prod_date = pd.to_datetime(raw_production["prod_date"]),  
14     units_total = lambda df: df["units_good"] + df["units_scrap"]  
15 )  
16  
17 stg_downtime = raw_downtime.assign(  
18     down_date = pd.to_datetime(raw_downtime["down_date"]),
```

# Dimensions : dim\_date, dim\_machine, dim\_product, dim\_shift

```
1 # dim_date : calendrier continu
2 dim_date = pd.DataFrame({"date": pd.date_range(stg_production["
    prod_date"].min(),
3                                     stg_production["
    prod_date"].max(),
4                                     freq="D")})
5 dim_date["date_sk"] = np.arange(1, len(dim_date)+1)
6 dim_date = dim_date.assign(
7     day=dim_date["date"].dt.day,
8     month=dim_date["date"].dt.month,
9     quarter=dim_date["date"].dt.quarter,
10    year=dim_date["date"].dt.year,
11    dow=dim_date["date"].dt.dayofweek
12 )
13 # dim_machine
```

## Tables de faits : fact\_production et fact\_downtime

```
1 # Keys de date
2 prod = stg_production.merge(dim_date[["date","date_sk"]],
3                               left_on="prod_date", right_on="date").
4                               drop(columns=["date"])
5 prod = (prod
6         .merge(dim_machine[["machine_id","machine_sk"]], on="machine_id")
7         .merge(dim_product[["sku","product_sk","ideal_rate_per_min"]], on="
8             sku")
9         .merge(dim_shift[["shift_id","shift_sk","planned_minutes"]], on="
10             shift_id"))
11
12 # Fait production (additive)
13 fact_production = prod[[
14     "date_sk","machine_sk","shift_sk","product_sk",
15     "units_good","units_scrap","units_total","ideal_rate_per_min",
16     "planned_minutes"
17 ]]
```



# Agrégations : Availability, Performance, Quality, OEE

```
1 # Agrégations par jour (toutes machines)
2 daily_prod = (fact_production
3     .merge(dim_date[["date_sk","year","month","date"]], on="date_sk")
4     .groupby(["date","year","month"], as_index=False)
5     .agg(units_good=("units_good","sum"),
6         units_total=("units_total","sum"),
7         planned_minutes=("planned_minutes","sum"),
8         ideal_rate=("ideal_rate_per_min","mean"))) # simplification
9
10 daily_down = (fact_downtime
11     .merge(dim_date[["date_sk","date"]], on="date_sk")
12     .groupby("date", as_index=False)["minutes_down"].sum())
13
14 daily = daily_prod.merge(daily_down, on="date", how="left").fillna({
15     "minutes_down":0})
16
17 # Indicateurs
18 daily["availability"] = (daily["planned_minutes"] - daily["
```

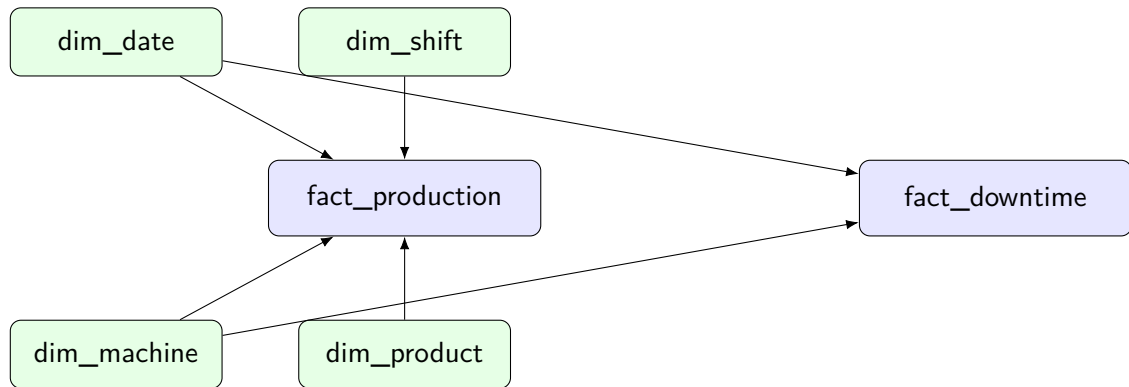
# Graphique : OEE mensuel (matplotlib)

```
1 x = range(len(oeo_month))
2 labels = [f"{int(y)}-{int(m):02d}" for y,m in zip(oeo_month["year"],
    oeo_month["month"])]
3 plt.figure(); plt.plot(x, oeo_month["oeo"], marker="o")
4 plt.title("OEE mensuel (exemple)")
5 plt.xlabel("Mois"); plt.ylabel("OEE")
6 plt.ylim(0,1); plt.xticks(x, labels, rotation=45, ha="right")
7 plt.tight_layout(); plt.savefig(out_dir/"monthly_oeo.png"); plt.show
    ()
8
9 # Persist
10 dim_date.to_csv(out_dir/"dim_date.csv", index=False)
11 dim_machine.to_csv(out_dir/"dim_machine.csv", index=False)
12 dim_product.to_csv(out_dir/"dim_product.csv", index=False)
13 dim_shift.to_csv(out_dir/"dim_shift.csv", index=False)
14 fact_production.to_csv(out_dir/"fact_production.csv", index=False)
15 fact_downtime.to_csv(out_dir/"fact_downtime.csv", index=False)
```

# TP : exercices proposés

- Calculer l'OEE par **machine** et par **zone** (area). Identifier les 5 machines les moins performantes.
- Décomposer l'OEE en **A/P/Q** par mois et tracer 3 courbes.
- Top 10 **modes de défaillance** par minutes perdues ; Pareto 80/20.
- Construire une **vue matérialisée** `v_oee_monthly(area)`.
- Ajouter des **règles DQ** : `minutes_down >= 0`, `units_total >= units_good`, `ideal_rate_per_min > 0`.

## Schéma en étoile — vue Production/Downtime



# Extension 1 — SCD2 sur dim\_machine

**Objectif** : historiser les changements (ex. area, status) d'une machine.

```
1 CREATE TABLE dim_machine (  
2     machine_sk          BIGINT PRIMARY KEY,  
3     machine_id          INT NOT NULL,  
4     machine_name        VARCHAR(40) NOT NULL,  
5     area                VARCHAR(40) NOT NULL,  
6     status              VARCHAR(20) DEFAULT 'Active',  
7     valid_from          DATE NOT NULL,  
8     valid_to            DATE NOT NULL,  
9     is_current          BOOLEAN NOT NULL,  
10    CONSTRAINT uk_machine_hist UNIQUE(machine_id, valid_from)  
11 );  
12  
13 -- R g l e : pour un changement, fermer l'ancienne ligne (valid_to =  
14 --           new_from - 1, is_current=false)  
15 -- et ins r e r la nouvelle (valid_from = date_changement, valid_to =  
16 --           '9999-12-31', is_current=true).  
17 -- Les faits pointent sur machine_sk via une table de mapping date
```

## Extension 2 — Accumulating Snapshot : fact\_workorder

**Objectif** : suivre le cycle de vie d'un ordre de maintenance (ouverture →clôture).

```
1 CREATE TABLE fact_workorder (  
2     wo_sk                BIGINT PRIMARY KEY,  
3     wo_number            VARCHAR(30) NOT NULL,    -- d g n r e  
4     machine_sk           BIGINT NOT NULL,  
5     open_date_sk         BIGINT,  
6     schedule_date_sk     BIGINT,  
7     start_date_sk        BIGINT,  
8     complete_date_sk     BIGINT,  
9     close_date_sk        BIGINT,  
10    labor_hours           NUMERIC(10,2),  
11    material_cost         NUMERIC(12,2),  
12    status               VARCHAR(20)  
13 );  
14 -- Colonnes de dates remplies au fil des vnements ; mesures mises  
    jour (upsert).
```

## Extension 3 — Bridge table : multi-causality des arrêts

**Objectif** : un arrêt peut avoir plusieurs causes. Utiliser un **bridge** avec pondération.

```
1 CREATE TABLE dim_failure_mode (  
2     failure_mode_sk BIGINT PRIMARY KEY,  
3     failure_mode     VARCHAR(60) UNIQUE  
4 );  
5  
6 CREATE TABLE br_downtime_cause (    -- bridge  
7     downtime_id      BIGINT NOT NULL,    -- identifiant du fait downtime  
8     failure_mode_sk BIGINT NOT NULL,  
9     allocation_pct   NUMERIC(5,4) NOT NULL CHECK (allocation_pct >= 0  
10    AND allocation_pct <= 1),  
11    PRIMARY KEY (downtime_id, failure_mode_sk)  
12 );  
13 -- Somme des allocation_pct = 1 par downtime_id.  
14 -- Les agrégations de minutes se font en r partissant minutes_down  
15    * allocation_pct.
```

## Extension 4 — Late Arriving Facts & Unknowns

- Créer des **membres inconnus** ( $SK = 0$  ou  $-1$ ) dans chaque dimension avec attributs "Unknown".
- Charger les faits en pointant vers  $SK=0$  si la dimension n'est pas encore disponible.
- **Rattrapage** : dès que la dimension arrive, *fixer* le SK dans le fait (update sécurisé).

```
1 -- Exemple membre inconnu
2 INSERT INTO dim_product(product_sk, sku, product_name,
   ideal_rate_per_min)
3 VALUES (0, 'UNK', 'Unknown Product', 1.0) ON CONFLICT DO NOTHING;
```



## Extension 5 — Data Quality : règles rapport d'anomalies

```
1  -- Exemples de règles
2  ALTER TABLE fact_downtime
3      ADD CONSTRAINT ck_minutes_nonneg CHECK (minutes_down >= 0);
4
5  -- Rapport d'anomalies (extrait)
6  CREATE VIEW dq_anomalies AS
7  SELECT 'NEGATIVE_MINUTES' AS rule, COUNT(*) AS n
8  FROM fact_downtime WHERE minutes_down < 0
9  UNION ALL
10 SELECT 'SCRAP_GT_TOTAL', COUNT(*)
11 FROM fact_production WHERE units_scrap > units_total
12 UNION ALL
13 SELECT 'ZERO_OR_NEG_RATE', COUNT(*)
14 FROM dim_product WHERE ideal_rate_per_min IS NULL OR
    ideal_rate_per_min <= 0;
```