

# Laporan Proyek Data Lakehouse - AdventureWorks

## 1. Manajemen Folder & File

Struktur folder proyek:

- data\_lake/adventureworks/files/ : menyimpan file sumber (csv, pdf, txt)
- data\_lake/adventureworks/tweets/ : menyimpan tweet mentah & hasil visualisasi
- database staging: warehouse\_temp\_sensor, market\_share\_report, external\_sentiment
- database warehouse: fact\_sales, dim\_customer, dim\_product, dim\_time, dim\_location

## 2. Desain Database

Database Staging:

- Tabel 'warehouse\_temp\_sensor': menyimpan data suhu dari gudang
- Tabel 'market\_share\_report': berisi hasil ekstraksi dari file PDF
- Tabel 'external\_sentiment': tweet tentang 'adventureworks'

Database Warehouse (Star Schema):

- fact\_sales: menyimpan metrik penjualan (misal dari file CSV gabungan)
- dim\_customer, dim\_product, dim\_time, dim\_location sebagai dimension tables

## 3. Jenis File & Analisis

- PDF: Ekstraksi teks laporan pasar (market share)
- CSV: Membaca data sensor suhu dan waktu
- TXT: Tweet teks diproses dan dibersihkan untuk wordcloud serta sentiment

## 4. Ringkasan Kode Utama

- data\_lake\_hans\_ingest.py: memproses file tweet dan menyimpan dalam berbagai format
- stagingdb.py: mengimpor semua file (csv, pdf, txt) ke dalam database staging
- analysis.py: melakukan analisis data
- tweet\_visualizer.py: menghasilkan visualisasi wordcloud

## 5. Diagram Star Schema

Laporan Proyek Data Lakehouse - AdventureWorks

